

Query lower bounds for log-concave sampling

Sinho Chewi* Jaume de Dios Pont† Jerry Li‡ Chen Lu§ Shyam Narayanan¶

October 31, 2023

Abstract

Log-concave sampling has witnessed remarkable algorithmic advances in recent years, but the corresponding problem of proving *lower bounds* for this task has remained elusive, with lower bounds previously known only in dimension one. In this work, we establish the following query lower bounds: (1) sampling from strongly log-concave and log-smooth distributions in dimension $d \geq 2$ requires $\Omega(\log \kappa)$ queries, which is sharp in any constant dimension, and (2) sampling from Gaussians in dimension d (hence also from general log-concave and log-smooth distributions in dimension d) requires $\tilde{\Omega}(\min(\sqrt{\kappa} \log d, d))$ queries, which is nearly sharp for the class of Gaussians. Here κ denotes the condition number of the target distribution. Our proofs rely upon (1) a multiscale construction inspired by work on the Kakeya conjecture in geometric measure theory, and (2) a novel reduction that demonstrates that block Krylov algorithms are optimal for this problem, as well as connections to lower bound techniques based on Wishart matrices developed in the matrix-vector query literature.

*School of Mathematics at Institute for Advanced Study, schewi@ias.edu. Part of this work was done while SC was a research intern at Microsoft Research.

†Department of Mathematics at University of California, Los Angeles, jdedios@math.ucla.edu.

‡Microsoft Research, jerrl@microsoft.com.

§Department of Mathematics at Massachusetts Institute of Technology, chenl819@mit.edu.

¶Department of Electrical Engineering and Computer Science at Massachusetts Institute of Technology, shyamsn@mit.edu. Part of this work was done while SN was a research intern at Microsoft Research.

Contents

1	Introduction	1
1.1	Our contributions	2
1.2	Related work	3
2	Technical overview	4
2.1	Geometric construction in low dimension	4
2.2	Lower bounds for sampling from Gaussians	8
2.2.1	Lower bound via Wishart matrices	8
2.2.2	Lower bounds via reduction to block Krylov	8
3	A general sampling lower bound in dimension two	11
3.1	Overview	11
3.2	Definitions and the information-theoretic argument	12
3.3	Reductions and properties of the construction	13
3.4	Construction of the distributions	17
4	A lower bound for sampling from Gaussians via Wishart matrices	19
4.1	Reducing inverse trace estimation to sampling	20
4.2	Lower bound for inverse trace estimation	21
4.3	Useful facts about Wishart matrices	23
5	A lower bound for sampling from Gaussians via reduction to block Krylov	23
5.1	Preliminaries	24
5.2	Lower bound against block Krylov algorithms	25
5.3	Reduction to block Krylov algorithms	28
5.3.1	Setup	28
5.3.2	Conditioning lemma	29
5.3.3	From query algorithms to block Krylov algorithms	32
A	Upper bound for log-concave sampling in constant dimension	41
B	Upper bound for sampling from Gaussians	43

1 Introduction

We study the problem of sampling from a target distribution on \mathbb{R}^d given query access to its unnormalized density. This is a fundamental algorithmic primitive arising in diverse fields, such as Bayesian inference, numerical simulation, and randomized algorithms [RC04]. Recently, there has been considerable progress in developing faster algorithms for this problem, particularly in the case where the target distribution is log-concave. In large part, these results have been achieved by exploiting the rich interplay between optimization and sampling [JKO98; Wib18], leading to novel sampling schemes inspired by classical optimization methods [Ber18; CLLMRS20; ZPFP20; LST21b; MCCFBJ21], as well as new quantitative convergence guarantees for sampling [Dal17; DMM19].

In light of such results, many prior works (e.g., [CCBJ18; LST21a; CBL22]) have raised the foundational question of whether the algorithmic upper bounds are tight. However, there is still a dearth of lower bounds for log-concave sampling. This lies in stark contrast to the analogous setting of convex optimization, in which the query complexity has been tightly characterized for a plethora of function classes [NY83; Nes18]. Such lower bounds yield important insights into the limitations of our existing algorithms and provide guidance towards identifying optimal ones.

Given the deep connections between the two fields, it is natural to ask why optimization lower bounds cannot be converted into sampling lower bounds. One way to do so is to directly reduce from optimization, as was done in [GLL22]. However, as we are interested in the intrinsic complexity of sampling, we make the standard assumption that the mode of the target distribution is zero to remove the optimization component of the sampling task, which rules out this approach. Another avenue is to borrow the techniques used for optimization lower bounds, but there are several obstructions to doing so. First, most optimization lower bounds hold against (classes of) deterministic algorithms and proceed by constructing specific adversarial functions [Bub15; Nes18]. In contrast, lower bounds for randomized algorithms are relatively recent and still not fully understood [WS17], which poses a major challenge for sampling algorithms, since they are inherently randomized. Second, whereas optimization constructions can employ local perturbations to hide the minima, sampling constructions need to hide the bulk of the mass of the target distribution, making them surprisingly delicate.

We now describe the problem in more detail. We consider the canonical setting in which target distribution π on \mathbb{R}^d is α -strongly log-concave and β -log-smooth, with its mode located at the origin. Namely, we assume $\pi \propto \exp(-V)$, where the potential $V : \mathbb{R}^d \rightarrow \mathbb{R}$ is twice continuously differentiable, α -strongly convex, β -smooth, and $\nabla V(0) = 0$. We let $\kappa := \beta/\alpha$ denote the *condition number* of π . We study algorithms in which the sampler is given query access to V and ∇V , and the goal is to produce a sample whose law is close to π in total variation distance. The complexity of the algorithm is measured by the number of queries made. Note that this oracle model captures the majority of sampling algorithms used in practice, including the unadjusted Langevin algorithm, Hamiltonian Monte Carlo, Metropolized random walks, and hit-and-run.

Despite the intense research activity centered on log-concave sampling, only a handful of works address the lower bound question, and the majority of them are either algorithm-specific or pertain to auxiliary problems such as estimation of the normalizing constant; see Section 1.2 for related work. To the best of our knowledge, currently the only general log-concave sampling lower bound is that of [CGLGR22], which establishes a sharp query lower bound of order $\Omega(\log \log \kappa)$ in dimension one. However, that work leaves open the question of obtaining stronger lower bounds in higher dimension, which is the more relevant case for applications. Even beyond the log-concave setting, we are aware of only one other work that obtains query lower bounds for sampling: the recent result of [CGLL23] is incomparable to the present work, as it considers a different setting, and we discuss it further in Section 1.2. Overall, the lack of sampling lower bounds points to a lack of tools for

addressing this problem and motivates the present work.

1.1 Our contributions

In this paper, we make significant progress on this problem by proving new lower bounds for sampling which reach beyond the one-dimensional setting considered in [CGLGR22]. In fact, for some settings of interest, our lower bounds match existing upper bounds up to constants, and we therefore obtain some of the first *tight* complexity results for sampling from log-concave distributions in dimension $d > 1$. We obtain lower bounds in two regimes:

Lower bounds in low dimension. Our first lower bound gives a tight characterization of the complexity of log-sampling in any constant dimension $d \geq 2$. We show:

Theorem 1 (informal, see Theorem 4). *For any dimension $d \geq 2$, any sampler for d -dimensional log-concave distributions with condition number κ requires $\Omega(\log \kappa)$ queries.*

Note that this result is exponentially stronger than the $\Omega(\log \log \kappa)$ lower bound in the univariate case [CGLGR22]. Moreover, when the dimension d is held fixed, we obtain a matching $O(\log \kappa)$ algorithmic upper bound, based on folklore ideas from the classical literature on sampling from convex bodies (Theorem 49). Together with the result of [CGLGR22] for $d = 1$, this settles the complexity of log-concave sampling in constant dimension.

On a technical level, the lower bound is based on a novel construction inspired by work on the Kakeya conjecture in geometric measure theory, which we believe may be of independent interest. We give a detailed description of the construction in Section 3.

Lower bounds in high dimension. Our second set of lower bounds applies to the high-dimensional setting and implies that when the dimension is sufficiently large, a polynomial dependence on the condition number κ is unavoidable (in contrast to Theorem 1, which only gives a logarithmic dependence on κ in low dimension). In fact, our lower bounds hold for the special case of sampling from Gaussians, for which they are nearly tight. We first prove the following theorem.

Theorem 2 (informal, see Corollary 19). *Any sampler for centered d -dimensional Gaussians with condition number κ requires $\Omega(\min(\sqrt{\kappa}, d))$ queries.*

We emphasize the fact that in our setting, the Gaussians are centered. Note that if the Gaussians were allowed to have varying means, then one can deduce a sampling lower bound by reducing the optimization task of minimizing a convex quadratic function $x \mapsto \langle (x - x_\star), \Sigma^{-1}(x - x_\star) \rangle$ to the task of sampling from the corresponding Gaussian $\mathcal{N}(x_\star, \Sigma)$. However, as previously alluded to, this does not address the inherent difficulty of the sampling problem.

The proof of Theorem 2 rests upon an elegant technique developed in the literature on the matrix-vector query model (see Section 1.2) in which the conditioning properties and sharp characterizations of the eigenvalue distribution of Wishart matrices are used to produce difficult lower bound instances for various tasks. We adapt this method to our context by reducing the task of inverse trace estimation to sampling (see Theorem 17).

As we show in Appendix B, the lower bound is nearly tight over the class of Gaussians, as it is possible to sample from a Gaussian using $O(\min(\sqrt{\kappa} \log d, d))$ queries using the block Krylov method. However, note that the lower bound from Theorem 2 does not match the block Krylov upper bound, and the lower bound of Theorem 2 is vacuous when κ is constant. In particular, it leaves open the possibility that the complexity of sampling from well-conditioned Gaussians is dimension-free. While such dimension-free rates are possible in convex optimization, our next result shows that the same is in fact not possible for log-concave sampling:

Theorem 3. (informal, see Theorem 45) *Let d be sufficiently large, and let $\kappa \leq d^{1/5-\delta}$. Then, any sampler for d -dimensional Gaussians with condition number κ requires $\Omega_\delta(\sqrt{\kappa} \log d)$ queries.*

In the regime for which Theorem 3 is valid, the lower bound matches the block Krylov upper bound up to constant factors, and hence we settle the complexity of sampling from Gaussians in this regime. Moreover, Theorems 2 and 3 together imply the *first* dimension-dependent lower bounds for general log-concave sampling. We conjecture that Theorem 3 holds for all κ for which $\sqrt{\kappa} \log d \leq d$, and we leave this question for future work.

Although Theorem 3 may appear to only be a mild improvement over Theorem 2, analyzing this regime is quite delicate, and we believe that the tools based on Wishart matrices employed in the proof of Theorem 2 may be insufficient to reach Theorem 3. Instead, we prove Theorem 3 by first establishing sharp lower bounds on the performance of block Krylov algorithms for the sampling task, and then providing a novel reduction (Lemma 39) which shows that block Krylov algorithms are optimal for this task. This reduction is quite general, and as the block Krylov algorithm and the matrix-vector query model are of wide interest in scientific computing and numerical linear algebra, we believe that our reduction may be broadly useful for tackling other problems in this space.

We remark that a concise way of summarizing Theorems 2 and 3 if we do not care about lower order terms is that sampling from Gaussians requires $\tilde{\Omega}(\min(\sqrt{\kappa} \log d, d))$ queries, where we write $f = \tilde{\Omega}(g)$ to mean $f = \Omega(g \log^{-O(1)}(g))$.

1.2 Related work

There is a vast literature on from sampling log-concave (and non-log-concave) distributions, and a full survey is beyond the scope of this paper. For a detailed exposition, see e.g. [Che22].

Lower bounds for log-concave sampling. As previously mentioned, the only unconditional lower bound against log-concave sampling is by [CGLGR22] for the one-dimensional setting, where the tight bound is $\Theta(\log \log \kappa)$. Other prior work on sampling lower bounds has fallen largely into one of several categories. One line of work studies lower bounds against a specific class of algorithm such as underdamped Langevin [CLW21] or MALA [CLACLR21; LST21a; WSC22]. However, these lower bounds techniques are tailored to the restricted class of algorithms that they consider and are not suitable for proving general query lower bounds. Another line of work considers lower bounds against computing normalizing constants [RV08; GLL20]. The work [Tal19] also investigates the computational complexity of sampling.

We mention two further lower bounds in different settings. The work of [CBL22] proves a lower bound against stochastic gradient oracles, and the work of [GLL22] proves a lower bound on the number of individual function value (i.e., zeroth-order) queries needed to sample from a density of the form $\exp(-\sum_{i \in I} f_i + \mu \|\cdot\|^2)$, where each f_i is convex, Lipschitz, and whose domain is the unit ball. In contrast, we consider deterministic, first-order oracle access. Moreover, their considerations are somewhat orthogonal to ours: [CBL22] focuses more on the role of noise, whereas we consider exact gradient access; and the lower bound of [GLL22] applies a direct reduction from optimization, which is also not in the spirit of the present work (in particular, we explicitly set the mode of the target distribution to zero).

Finally, we also mention the recent work [CGLL23], which proves query lower bounds for non-log-concave sampling in a different metric (the Fisher information). This work is inspired by the corresponding upper bounds of [BCESZ22] and can be viewed as lower bounds against *local mixing*.

Upper bounds for log-concave sampling. Starting with the seminal papers of [DT12; Dal17; DM17], there has been a flurry of recent work on proving non-asymptotic guarantees for log-concave sampling, with iteration complexities that scale polynomially in the condition number and

dimension. This includes analyses for the classical Langevin dynamics [Wib18; DK19; DMM19; VW19; BCESZ22; CELSZ22; AT23], mirror and proximal methods [Wib19; CLLMRS20; SR20; ZPFP20; AC21; Jia21; LST21b; CCSW22; CE22; GV22; LTVW22; FYC23], the Metropolis-adjusted Langevin algorithm (MALA) [DCWY18; CDWY20; LST20; CLACLR21; WSC22; AC23], and many others [CCBJ18; SL19; DR20; DLLW21; MCCFBJ21].

Our upper bound for sampling from Gaussians (Theorem 52) is closely related to the use of the conjugate gradient algorithm for sampling from Gaussians [NS22]. Also, our $O(\log \kappa)$ upper bound algorithm is closely related to rounding procedures which have been previously used in the convex body sampling literature (see, e.g., [LV06]).

Matrix-vector product query model. While matrix-vector queries have been studied in scientific computing for decades (e.g., [BFG96]), they have only been studied in the theoretical computer science literature recently, with a fully formalized model described in [SWYZ19]. The most relevant works to ours are those that study the matrix-vector query complexity of spectral properties, such as estimating top eigenvectors [SAR18; BHSW20], trace and matrix norms [Hut90; WWZ14; RWZ20; DM21; MMMW21], the full eigenspectrum [CKSV18; BKM22], and low-rank approximation [MM15; BCW22]. We remark that the non-adaptive matrix-vector product model is closely related to sketching, which has enjoyed a large body of work (see, e.g., [Woo14] for a survey).

2 Technical overview

Here we summarize the main technical ideas used to prove our lower bounds. For details, see Section 3 for Theorem 1, Section 4 for Theorem 2, and Section 5 for Theorem 3.

2.1 Geometric construction in low dimension

Theorem 1 is proved with a construction in dimension two. For convenience, in this section we use radial coordinates to denote points in \mathbb{R}^2 , so $\omega := (x, y) = (r, \theta)$, where $r \in \mathbb{R}_+$ and $\theta \in [0, 2\pi)$. We denote sectors of \mathbb{R}^2 enclosed by angles θ_1 and θ_2 as $S(\theta_1, \theta_2) := \{(r, \theta) \in \mathbb{R}^2 : \theta \in [\theta_1, \theta_2]\}$, and denote bounded sectors as $S_{\text{bdd}}(\theta_1, \theta_2, r) := \{(r', \theta) \in \mathbb{R}^2 : \theta \in [\theta_1, \theta_2], r' \leq r\}$.

The argument is information-theoretic in nature. We will construct a family of strongly log-concave and log-smooth distributions $\{\pi_1, \dots, \pi_m\}$, where each $\pi_b \propto \exp(-V_b)$, which satisfies two key properties. First, different distributions π_b and $\pi_{b'}$ are well separated in total variation distance; and second, if b is chosen uniformly at random from $[m]$, then querying the potential $(V_b(\omega), \nabla V_b(\omega))$ at any $\omega \in \mathbb{R}^2$ will reveal $O(1)$ bits of information about b . The lower bound in Theorem 1 follows readily from the existence of such a family, provided that m and κ are polynomially related. On the one hand, because the distributions are well-separated in total variation, if we can sample well from the distribution π_b using queries, we can identify the index b with high probability. On the other hand, because there are m distributions and every query reveals $O(1)$ bits of information about b , we need at least $\Omega(\log m) = \Omega(\log \kappa)$ queries to identify b , which results in a $\Omega(\log \kappa)$ query lower bound for log-concave sampling.

How do we construct such a family? A first attempt is to consider distributions supported on thin convex sets that have no overlap. For $b = \frac{1}{\kappa}, \frac{2}{\kappa}, \dots, 1$, let $\pi_b = \text{uniform}(\mathcal{Z}_b)$, where $\mathcal{Z}_b = S_{\text{bdd}}(\frac{\pi}{2}b, \frac{\pi}{2}(b + \frac{1}{2\kappa}), 1)$, and the size of the family is $m = \lfloor \kappa \rfloor$. The potential V_b is the convex indicator of \mathcal{Z}_b , i.e., it is 0 on \mathcal{Z}_b and $+\infty$ outside. Morally, the distributions π_b can be thought of as having condition number κ .

This family does satisfy the two properties needed for the lower bound: different distributions are certainly well-separated because they have disjoint supports; and when we query any potential

V_b at a point $\omega \in \mathbb{R}^2$, we always receive one bit of information: whether or not ω lies in the support of π_b . However, the distributions in this family are neither strongly log-concave nor log-smooth. It is easy to make them strongly log-concave while still satisfying the desired properties: we can adjust the distributions by adding the same quadratic function $\frac{\|\cdot\|^2}{2}$ to all of the potentials V_b . But it is much harder to make this family log-smooth.

One way to make this construction log-smooth is to let the potentials V_b grow slowly (linearly) to infinity outside of the their zero sets \mathcal{Z}_b , which leads to a modified second attempt: for $m = \kappa^{\Omega(1)}$, $b = \frac{1}{m}, \dots, 1$, let π_b have potential $V_b = \tilde{V}_b + \frac{\|\cdot\|^2}{2\kappa^{O(1)}}$, where $\mathcal{Z}_b = S(\frac{\pi}{2}b, \frac{\pi}{2}(b + \frac{1}{2m}))$, and $\tilde{V}_b(\omega) = \kappa \text{dist}(\omega, \mathcal{Z}_b)$. Note that the potentials V_b are in fact still not smooth at the boundaries of the sets \mathcal{Z}_b , but this can be fixed by mollifying V_b . The distributions in this family will be well-separated, because an $\Omega(1)$ fraction of the mass of π_b will lie in \mathcal{Z}_b , and the sets \mathcal{Z}_b are disjoint for different b . Unfortunately, this family no longer reveals $O(1)$ bits per query: for any $\omega \in \mathbb{R}^2$, we can identify b with a single query to $(V_b(\omega), \nabla V_b(\omega))$, because either $\omega \in \mathcal{Z}_b$, or $\nabla V_b(\omega)$ reveals the direction of \mathcal{Z}_b , and in both cases the index b itself is identified.

We can reduce the information revealed by queries by more carefully controlling the growth of \tilde{V}_b , so that the further away a point ω lies from \mathcal{Z}_b , the fewer the number of bits will be revealed by $(\tilde{V}_b(\omega), \nabla \tilde{V}_b(\omega))$. This motivates a third attempt at the construction. For $m = 2^N = \kappa^{\Omega(1)}$, $b = \frac{1}{m}, \dots, 1 - \frac{1}{m}$, let $b = 0.b_1 \dots b_N$ be the binary expansion of b , and let $[b]_k = 0.b_1 \dots b_k$ be the truncation of b up to the k -th bit. For $k = 1, \dots, N$, let $\mathcal{Z}_{k,b}^{\text{radial}} = S(\frac{\pi}{2}[b]_k, \frac{\pi}{2}([b]_k + 2^{-k}))$, and let $\phi_{k,b}^{\text{radial}}(x) = \kappa^{O(1)} 2^{-k} \text{dist}(x, \mathcal{Z}_{k,b}^{\text{radial}})$. Finally, let $V_b^{\text{radial}} = \frac{\|\cdot\|^2}{2\kappa^{O(1)}} + \tilde{V}_b^{\text{radial}}$, where

$$\tilde{V}_b^{\text{radial}} = \max_{k=1, \dots, N} \phi_{k,b}^{\text{radial}}.$$

The potentials V_b^{radial} will again have to be mollified to be made smooth. It turns out that the potentials $\tilde{V}_b^{\text{radial}}$ will grow fast enough outside $\mathcal{Z}_{N,b}^{\text{radial}}$ such that the distributions will be well-separated. It also turns out that queries indeed reveal $O(1)$ bits of information on average. This can be seen as follows: note that the sets $\mathcal{Z}_{k,b}^{\text{radial}}$ are sectors such that $\mathcal{Z}_{k,b}^{\text{radial}} \supset \mathcal{Z}_{k+1,b}^{\text{radial}}$, and as k increases, $\mathcal{Z}_{k,b}^{\text{radial}}$ becomes thinner around the ray $\{\theta = \frac{\pi}{2}b\}$; also note that as k increases, the growth rate of $\phi_{k,b}^{\text{radial}}$ outside its zero set $\mathcal{Z}_{k,b}^{\text{radial}}$ is decreasing; these two properties imply that if we query a point $\omega = (r, \theta)$ that is far from the sector $\mathcal{Z}_{i,b}^{\text{radial}}$ (in the sense that $\theta \notin [\frac{\pi}{2}[b]_i - 100 \cdot 2^{-i}, \frac{\pi}{2}[b]_i + 100 \cdot 2^{-i}]$), then the value of $\tilde{V}_b^{\text{radial}}(\omega)$ will not depend on any $\phi_{k,b}^{\text{radial}}$ for $k > i$, and hence querying $\tilde{V}_b^{\text{radial}}(\omega)$ will only reveal b up to the i -th bit. As a result, if b is chosen uniformly, then for a fixed query ω with high probability we will have $\omega \notin \mathcal{Z}_{k,b}^{\text{radial}}$ for any $k = O(1)$, so the query will only reveal $O(1)$ bits of information about b .

Yet this construction fails because of the mollification step, which we have so far ignored. To make the potentials V_b smooth, we will instead take $V_b = \chi_\delta * \tilde{V}_b^{\text{radial}} + \frac{\|\cdot\|^2}{2\kappa^{O(1)}}$, where χ_δ is supported on a ball of radius $\delta < 2^{-2N}$. We would hope that the potential $\chi_\delta * \tilde{V}_b^{\text{radial}}$ still satisfies the property that querying a point $\omega = (r, \theta)$ that is far from $\mathcal{Z}_{i,b}^{\text{radial}}$ only reveals b up to the i -th bit. When r is not too close to the origin (say $r > 100 \cdot 2^{-i}$), this is indeed still true: if ω satisfies $\theta \notin [\frac{\pi}{2}[b]_i - 200 \cdot 2^{-i}, \frac{\pi}{2}[b]_i + 200 \cdot 2^{-i}]$, then the entire δ -neighbourhood of ω will satisfy $\theta \notin [\frac{\pi}{2}[b]_i - 100 \cdot 2^{-i}, \frac{\pi}{2}[b]_i + 100 \cdot 2^{-i}]$, so the value of $\tilde{V}_b^{\text{radial}}$ on the δ -neighbourhood of ω will not depend on any $\phi_{k,b}^{\text{radial}}$ for $k > i$, hence the value of $(\chi_\delta * \tilde{V}_b^{\text{radial}})(\omega)$ will also not reveal any information of b beyond the i -th bit. But when ω is very close to the origin ($r < \delta$), the δ -neighbourhood of ω will intersect $\mathcal{Z}_{N,b}^{\text{radial}}$, which means that the value of $(\chi_\delta * \tilde{V}_b^{\text{radial}})(\omega)$ will depend on $\phi_{k,b}^{\text{radial}}$ for all k and hence on all bits of b . In other words, mollification leaks information around the origin. As a result, if we query points δ -close to the origin, we will again identify b in a single query.

The way to resolve the leakage at the origin is to create a branching structure, such that all V_b are equal near the origin so that no information is leaked at small scales, and such that far away from the origin V_b is small around the ray $\{\theta = \frac{\pi}{2}b\}$ so that π_b still concentrates around different sectors. We keep the choices of m and b from the previous construction. The potentials will be $V_b = \chi_\delta * \tilde{V}_b + \frac{\|\cdot\|^2}{2\kappa^{O(1)}}$, where $\tilde{V}_b = \max_{k=1,\dots,N} \phi_{k,b}$, and $\phi_{k,b}(\omega) = \kappa^{O(1)} 2^{-k} \text{dist}(\omega, \mathcal{Z}_{k,b})$. The zero set $\mathcal{Z}_{k,b}$, instead of being a radial sector like $\mathcal{Z}_{k,b}^{\text{radial}}$, is now thickened adaptively.

We intuitively describe how to generate $\mathcal{Z}_{k,b}$. Each $\mathcal{Z}_{k,b}$ will be a thickening of $\mathcal{Z}_{k,b}^{\text{radial}}$, by simply including all points within some distance d_k of $\mathcal{Z}_{k,b}^{\text{radial}}$. We define $\mathcal{Z}_{\leq k,b} := \bigcap_{k' \leq k} \mathcal{Z}_{k',b}$: note that each $\mathcal{Z}_{\leq k,b}$ is getting smaller as k increases, and $\mathcal{Z}_{\leq N,b}$ is the zero set of \tilde{V}_b .

Consider some radii $r_0 < r_1 < r_2 < \dots$. To generate $\mathcal{Z}_{1,b}$, we thicken $\mathcal{Z}_{1,b}^{\text{radial}}$ (corresponding to the radial sector matching on the first bit), so that it contains $S_{\text{bdd}}(0, \pi/2, r_0)$ (corresponding to the quarter-circle near the origin). This avoids leaking information near the origin, as every x within radius r will be in $\mathcal{Z}_{1,b}$, which means $\phi_{1,b}$ will also be 0. Indeed, we can thicken $\mathcal{Z}_{1,b}^{\text{radial}}$ just the right amount so that it contains $S_{\text{bdd}}(0, \pi/2, r_0)$. For the concrete example where $N = 4$, and $b = 0.1010$, we show a description of $\mathcal{Z}_{1,b}$ in Figure 1a: we shade $S_{\text{bdd}}(0, \pi/2, r_0)$ in dark blue, $\mathcal{Z}_{1,b}^{\text{radial}} = S(\pi/4, \pi/2)$ in medium blue, and the additional thickening required in light blue.

To generate $\mathcal{Z}_{k,b}$ for $k \geq 2$, we thicken a much thinner angular sector. This ensures that at large radii, the arc of $\mathcal{Z}_{k,b}$ is not too big. We will inductively thicken $\mathcal{Z}_{k,b}$ by some amount d_k just enough to contain $\mathcal{Z}_{k-1,b} \cap S_{\text{bdd}}(0, \pi/2, r_{k-1})$. Consider one more example for $k = 2$ (again for $N = 4$, and $b = 0.1010$), in Figure 1b. Note that $\mathcal{Z}_{2,b}^{\text{radial}}$ is the sector $S(\frac{\pi}{4}, \frac{3\pi}{8})$ (shaded in medium blue), and the thickened region (in light blue emanating from both sides of the sector) is just enough to capture all of $\mathcal{Z}_{1,b}$ that was within radius r_1 . However, for larger radii, $\mathcal{Z}_{2,b}$ is much thinner than $\mathcal{Z}_{1,b}$. In addition, if we know the first bit $b_1 = 1$, then querying V_b anywhere in $\{r \leq r_1\}$ will not reveal any information about the second bit b_2 . This is because either we were in $\mathcal{Z}_{1,b}$ which only depends on b_1 (in which case $\phi_{1,b} = \phi_{2,b} = 0$ as we thickened to make sure $\mathcal{Z}_{2,b} \supset \mathcal{Z}_{1,b} \cap S_{\text{bdd}}(0, \pi/2, r_1)$), or we weren't, in which case $\phi_{1,b}$ grows much more quickly than $\phi_{2,b}$.

We can also continue this process inductively for $k = 3, 4$ (Figures 1c and 1d): we show $\mathcal{Z}_{\leq k,b}$. The intuition for why this prevents leaking of information near the origin is that even if k is large, $\mathcal{Z}_{k,b}$ in the smaller-radius regions is decided by $\mathcal{Z}_{k',b}$ for $k' \ll k$, so we cannot learn any later bits.

The comparisons of $\mathcal{Z}_{k,b}^{\text{radial}}$ and $\mathcal{Z}_{\leq k,b}$ for $b = 0.1010$ and for all $k \leq 4$ are shown together in Figure 1. The picture is not to scale, and the radial arcs represent the radii $r_i = 2^i r_0$, for $i = 0, \dots, 4$.

The construction of $\mathcal{Z}_{\leq k,b}$ means that for $k > 1$, querying $\phi_{k,b}$ within $\{r \leq 2^{k-1} r_0\}$ will not reveal the k -th bit, and so even querying the mollified $\chi_\delta * \phi_{k,b}$ within $\{r \leq 2^{k-2} r_0\}$ will not reveal the k -th bit, which stops information leaking near the origin.

Since $\tilde{V}_b = \max_{k=1,\dots,N} \phi_{k,b}$, the zero set of \tilde{V}_b coincides with $\mathcal{Z}_{\leq N,b}$, and for the choice of $b = 0.1010$, this is shown in the first panel of Figure 2. It turns out that each π_b will concentrate around the zero set of \tilde{V}_b , and the other panels of Figure 2 show these zero sets for seven different values of b in the set $\{\frac{1}{16}, \dots, \frac{15}{16}\}$ at larger scales. We can see that far out from the origin the zero sets become well-separated, and hence the distributions are well-separated in total variation.

We already discussed how the thickening of $\mathcal{Z}_{k,b}$ means that querying $\phi_{k,b}$, and hence \tilde{V}_b , near the origin will not reveal the higher bits of b . For query points $\omega = (r, \theta)$ where r is large, the same analysis on $\tilde{V}_b^{\text{radial}}$ tells us that $\tilde{V}_b(x)$ (even after mollification) will reveal $O(1)$ bits of information about b when b is chosen uniformly. As mentioned earlier, such a family of distributions readily leads to a sampling lower bound of $\Omega(\log m)$, where m is the size of the family. Since we can choose $m = \kappa^{\Omega(1)}$, this leads to the $\Omega(\log \kappa)$ lower bound. Details of the proof can be found in Section 3.

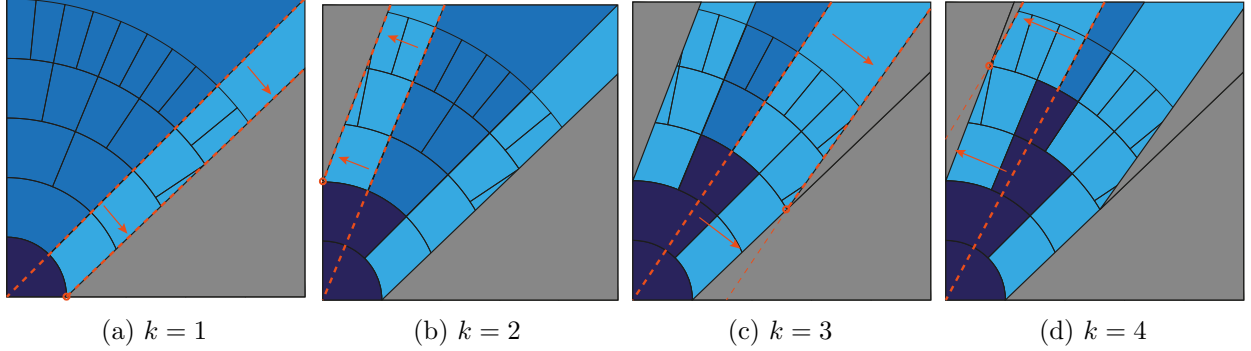


Figure 1: Comparison of $Z_{\leq k,b}^{\text{radial}}$ (the sector in medium blue) with $Z_{k,b}$ (union of dark, medium, and light blue), for $k = 1, 2, 3, 4$, and $b = 0.1010$. Dark blue represents the larger angular sectors closer to the origin, and light blue represents the additional fattening from taking sumsets. Each $Z_{k,b}$ is constructed by thickening $Z_{k,b}^{\text{radial}}$ enough (illustrated by the red arrows) such that no information about the k -th bit is revealed close to the origin, but $Z_{k,b}$ continues to get thinner at large radii.

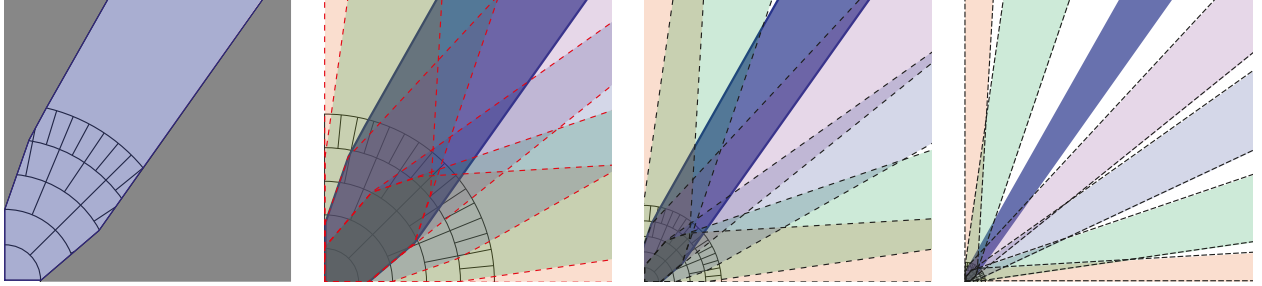


Figure 2: Zeros sets of \tilde{V}_b . The first panel shows the zero set for $b = 0.1010$. The other panels show the zeros sets for different values of b at different scales. Note that far away from the origin the zero sets become well-separated, which leads to the distributions being well-separated in total variation. Note that if b, b' match in the first ℓ bits, then they will agree up to the ℓ -th circle, as those circles only depend on $Z_{\leq \ell, b}$ even for ℓ much less than K .

Connections to Kakeya constructions. The construction outlined above is related to Perron’s construction [Per28] of Besicovich (Kakeya) sets known as *Perron trees*. Kakeya sets are sets with area zero that contain the translation of a unit segment in any direction. While Kakeya sets over finite fields have been investigated before in theoretical computer science, e.g., [SS08; Dvi09; Juk11], our construction is inspired by Kakeya sets over continuous domains, namely \mathbb{R}^2 . To our knowledge, this is one of the first applications of these geometric ideas to theoretical computer science.

There are many similarities between our construction and that of Perron. Perron’s construction proceeds by the method of sprouting. Sprouting is an iterative process in which, at each step, one adds further and further smaller triangles to the pre-existing construction. The figure is then rescaled in order to have height 1. The construction after n steps contains 2^n triangles of small aperture $\Omega(2^{-n})$, and has area $O(n^{-1})$. We do a similar process in the definition of our sets $Z_{k,b}$, and indeed, ultimately our hard instance has a very similar tree-like structure.

While we were inspired by the construction of Perron trees, there are also key differences between our hard instance and Perron’s construction. Indeed, in our setting, we need to minimize overlap (so that the resulting distributions are well-separated) while simultaneously ensuring that information is not leaked by queries. In contrast, Kakeya sets are explicitly designed to maximize overlap. Secondly, the iterates of Perron trees are convex sets, not convex functions. One must turn these

convex sets into convex functions somehow. This is additionally complicated by the fact that these iterates are not nested. In our construction, we must take great care to create nested convex sets, so that the resulting functions are convex and still maintain the structure of the sets.

2.2 Lower bounds for sampling from Gaussians

We now turn to our lower bounds against sampling from Gaussians. Recall that our goal is to provide a lower bound on sampling from a Gaussian $\mathcal{N}(0, \Sigma)$, where $\Lambda := \Sigma^{-1}$ has condition number κ . Note that the corresponding potential is $V(x) = \frac{1}{2} \langle x, \Lambda x \rangle$, and we are allowed zeroth-order and first-order queries, which means for a query x , we receive $x^\top \Lambda x$ and Λx . Hence, adaptive queries are equivalent to adaptive matrix-vector product computations with Λ .

The first observation we make is that we can reduce the problem of sampling from the Gaussian to estimating the trace of Σ . This is because if X is a sample from a distribution which is close in total variation distance to $\mathcal{N}(0, \Sigma)$, then $\|X\|_2^2 \approx \text{tr}(\Sigma)$ with high probability. Therefore, it suffices to demonstrate a lower bound for the following problem: given matrix-vector product computations with Λ , approximately compute $\text{tr}(\Lambda^{-1})$.

2.2.1 Lower bound via Wishart matrices

For any d , let $W \in \mathbb{R}^{d \times d}$ have the Wishart(d) distribution. That is, $W = XX^\top$, where $X \in \mathbb{R}^{d \times d}$ has i.i.d. $\mathcal{N}(0, 1/d)$ entries. We take W to be the precision matrix, $\Lambda = W$. Our first lower bound shows that $\Omega(d)$ matrix-vector queries with W are necessary to estimate the trace of W^{-1} even to constant multiplicative accuracy, with constant success probability (Theorem 18). Since the condition number of W is $\Theta(d^2)$ with high probability, we obtain one extreme of the claimed lower bound $\Omega(\min(\sqrt{\kappa}, d))$. The general lower bound for all κ then follows from a padding argument.

This lower bound approach is inspired by [BHSW20], which proved a query lower bound for estimating the minimum eigenvalue of W . Their approach relies on the fact that if we condition on any sequence of $(1 - \Omega(1))d$ adaptive queries, the posterior distribution of the remaining eigenvalues behaves similarly to the original distribution of the eigenvalues of W . In addition, while the smallest eigenvalue of W is usually about $1/d^2$, its distribution has heavy tails: with probability $\Theta(\sqrt{\varepsilon})$, the smallest eigenvalue of W is below ε/d^2 . Consequently, even conditioned on $d/2$ adaptive queries, we are unable to learn the minimum eigenvalue up to a constant factor with high probability.

In our setting, we instead wish to show that learning the trace of W^{-1} is hard. However, the smallest eigenvalue of the Wishart matrix is so small that with high probability, $\text{tr}(W^{-1}) = \Theta(\lambda_{\min}(W)^{-1})$. While most of the time the trace is $O(d^2)$, with probability $\Theta(\sqrt{\varepsilon})$ the posterior distribution of the smallest eigenvalue of W after our adaptive queries may be ε/d^2 . Hence, we will be unable to determine whether the trace is $\leq O(d^2)$ or $\geq \Omega(d^2/\varepsilon)$ with high probability.

This lower bound technique is clean and nearly optimal, but as previously mentioned it is vacuous (of constant order) when $\kappa = O(1)$, whereas we expect the complexity of the problem to increase as $d \rightarrow \infty$. To tackle this setting, we introduce a second approach.

2.2.2 Lower bounds via reduction to block Krylov

Our second technique works in two parts. First, we show that for a specific hard distribution over instances, any block Krylov-style algorithm requires $\Omega(\min(\sqrt{\kappa} \log d, d))$ queries to estimate $\text{tr}(\Sigma)$. Then, we show a general purpose reduction which demonstrates that for this hard instance (and indeed, any rotationally invariant instance), block Krylov methods are actually optimal.

Lower bound for block Krylov algorithms. Recall the block Krylov technique: the algorithm chooses K i.i.d. random vectors $v_1, \dots, v_K \sim \mathcal{N}(0, I)$, and computes $\Lambda^j v_k$ for all $j \leq T, k \leq K$. This

can be done using KT adaptive queries, by querying $\Lambda^j v_k$ to learn $\Lambda^{j+1} v_k$. For our purposes, it suffices to consider block Krylov algorithms with $K = T$ and to prove a lower bound on the smallest number K needed to successfully estimate $\text{tr}(\Sigma)$, for $\Sigma = \Lambda^{-1}$.

We will construct two diagonal matrices D, D' with all eigenvalues between 1 and κ , such that $\text{tr}(D^{-1})$ and $\text{tr}((D')^{-1})$ are sufficiently different. In addition, if Λ, Λ' are random rotations of D, D' , respectively, then $\{\Lambda^j v_k\}_{j,k \leq K}$ and $\{(\Lambda')^j v_k\}_{j,k \leq K}$ are hard to distinguish for $K \leq c\sqrt{\kappa} \log d$ for a small constant c (Lemma 32). Thus, unless $K \geq \Omega(\sqrt{\kappa} \log d)$, we cannot estimate the trace.

To explain the intuition behind Lemma 32, we first consider what happens if we only have $\{\Lambda^j v\}_{j \leq K}$ for a single random vector v (i.e., power method). Letting $\lambda_1, \dots, \lambda_d$ be the eigenvalues of Λ , we have $\Lambda^j v = \sum_{i=1}^d \lambda_i^j \alpha_i u_i$, where u_i is the i -th eigenvector of Λ and $v = \sum_{i=1}^d \alpha_i u_i$. Intuitively, the only information we obtain from these vectors are their pairwise inner products, since we could have randomly rotated Λ . Therefore, the only information we have is $\langle \Lambda^j v, \Lambda^{j'} v \rangle = \sum_{i=1}^d \lambda_i^{j+j'} \alpha_i^2$, which is the set $\{\sum_{i=1}^d \lambda_i^j \alpha_i^2\}_{j \leq 2K}$. Since v is random, we may think of all of the α_i^2 as 1 for simplicity, and so we know $\{\sum_{i=1}^d \lambda_i^j\}_{j \leq 2K}$. Our goal is to use this information to learn $\text{tr}(\Lambda^{-1}) = \sum_{i=1}^d \lambda_i^{-1}$.

We connect this to the problem of estimating $1/x$ as a linear combination of $1, x, x^2, \dots, x^K$, a classic problem in approximation theory that is often tackled with *Chebyshev polynomials*. Indeed, this relation to Chebyshev polynomials is the main tool in the analysis of essentially all Krylov methods. In our setting, as we desire lower bounds, we apply the fact that Chebyshev polynomials are *optimal* in generating certain approximations. More concretely, suppose that there are only K distinct eigenvalues $\lambda_1, \dots, \lambda_K$, with each λ_i having some multiplicity N_i . Since we want to show that estimating $\text{tr}(\Lambda^{-1})$ is hard, this amounts to showing that knowing $\sum_{i=1}^K N_i \lambda_i^j$ for $0 \leq j \leq K$ is insufficient to learn $\sum_{i=1}^K N_i / \lambda_i$. We express this as a linear program (if we relax the N_i to be reals), the dual of which precisely captures whether $1/x$ can be approximated well by a degree- K polynomial at $\lambda_1, \dots, \lambda_K$ (Proposition 29). If we choose the λ_i to be the local extrema of a degree- K Chebyshev polynomial, shifted so that $\lambda_1 = 1$ and $\lambda_K = \kappa$, then it is known that one cannot estimate $1/x$ up to error $d^{-\Omega(1)}$ at these points (which is needed for trace estimation), unless $K \geq \Omega(\sqrt{\kappa} \log d)$. At a high level, this is the reason why we need $\Omega(\sqrt{\kappa} \log d)$ iterations of the power method.

For general block Krylov algorithms, the algorithm obtains $\langle v_\ell, \Lambda^j v_k \rangle$, for $0 \leq j \leq K$ and $1 \leq k, \ell \leq K$. Now, the information that the algorithm sees is captured by the matrices $\{\langle v_\ell, \Lambda^j v_k \rangle\}_{k, \ell \leq K}$, for $j = 1, \dots, K$. Here, we show that provided K is sufficiently small compared to d , we can find choices of multiplicities N_1, \dots, N_K and N'_1, \dots, N'_K , such that the corresponding matrices D, D' have significantly different traces (i.e., $\sum_{i=1}^K (N_i - N'_i) / \lambda_i$ is large) but the information from queries is not enough to distinguish between Λ and Λ' , which we establish via a coupling argument.

Reduction to block Krylov algorithms. The argument outlined above shows block Krylov algorithms with $K = o(\sqrt{\kappa} \log d)$ cannot distinguish between two families of randomly rotated matrices with difference traces (Λ coming from D and Λ' coming from D'), and hence cannot solve the trace estimation task. Our next technical contribution is a reduction which allows us to simulate the output of any *adaptive* algorithm with K queries on our hard instance, given only the responses to a block Krylov algorithm. Thus, a lower bound against block Krylov methods translates into a lower bound against any query algorithm. We now give a high-level description of the reduction.

Since we prove lower bounds based on randomized constructions, it suffices to consider adaptive deterministic algorithms, i.e., each query v_k is a deterministic function of the previous queries and oracle outputs. The difficulty of proving such a lower bound against such an algorithm is the adaptivity of the queries, which makes it difficult to reason about how much information the algorithm has learned. However, since our lower bound construction for block Krylov algorithms is rotationally invariant, intuitively the adaptivity does not help: the algorithm may as well query a random direction which it has not yet explored.

However, this intuition is not entirely correct: if the algorithm has previously queried a vector v and received the information Λv , then it may be useful to query Λv in order to receive the information $\Lambda^2 v$, instead of querying a completely random new direction. Indeed, computing powers $v, \Lambda v, \Lambda^2 v, \dots$ is precisely the essence of the power method, as discussed above. To account for this, we move to the following stronger oracle model: if the algorithm has selected vectors v_1, \dots, v_k , then at iteration k it receives all of the information $(\Lambda^i v_j)_{i+j \leq k}$ for free. Now, there is provably no benefit to querying vectors which lie in the span of the previous queries and oracle outputs.

Recall that our goal is to argue that an adaptive deterministic algorithm can be simulated by an algorithm which simply makes i.i.d. Gaussian queries z_1, z_2, \dots, z_K , in the following sense. In the stronger oracle model, at iteration k , the adaptive algorithm has made queries $(v_1^{\text{alg}}, \dots, v_k^{\text{alg}})$ and received information $(\Lambda^i v_j^{\text{alg}})_{i+j \leq k}$ and it picks a new vector v_{k+1} which lies orthogonal to its received information. Suppose that using only the Gaussian queries z_1, z_2, \dots, z_k , we have simulated queries $v_1^{\text{sim}}, v_2^{\text{sim}}, \dots, v_k^{\text{sim}}$ which are equivalent to the execution of the adaptive algorithm in the sense that the law of the information $(\Lambda^i v_j^{\text{sim}})_{i+j \leq k}$ is precisely the same as the law of the algorithm's information $(\Lambda^i v_j^{\text{alg}})_{i+j \leq k}$. Since the algorithm is deterministic, v_k^{alg} is a function $v_k((\Lambda^i v_j^{\text{alg}})_{i+j < k})$ of algorithm's accumulated information. Thus, in order to simulate the adaptive algorithm for one more step, it is natural to consider taking $v_k^{\text{sim}} := v_k((\Lambda^i v_j^{\text{sim}})_{i+j < k})$. However, we will be unable to compute $\Lambda^i v_k^{\text{sim}}$ for any $i \geq 1$, because the simulation must be based on the Gaussian queries z_1, z_2, \dots, z_k , whereas this definition of v_k^{sim} requires making queries at $v_1^{\text{sim}}, v_2^{\text{sim}}, \dots, v_{k-1}^{\text{sim}}$.

Thus far, we have not invoked the rotational invariance of Λ , which is crucial to the argument. The key is that although we cannot directly take $v_k((\Lambda^i v_j^{\text{sim}})_{i+j < k})$ to be our next simulated point, we can *rotate* \tilde{v}_k into $v_k((\Lambda^i v_j^{\text{sim}})_{i+j < k})$ via a unitary matrix U_k ; moreover, we can arrange that U_k fixes all of the previous information $(\Lambda^i v_j^{\text{sim}})_{i+j < k}$, because $v_k((\Lambda^i v_j^{\text{sim}})_{i+j < k})$ lies orthogonal to this information (recall, we can assume that each deterministic function $v_k(\cdot)$ outputs a vector orthogonal to its inputs, due to our choice of oracle model). The intuition is that due to the rotational invariance of Λ , then conditioned on the data $(\Lambda^i v_j^{\text{sim}})_{i+j < k}$, the distribution of Λ is still rotationally invariant on the orthogonal subspace of the data; hence, $U_k \tilde{v}_k = v_k((\Lambda^i v_j^{\text{sim}})_{i+j < k})$ ought to have the same law as \tilde{v}_k , i.e., querying the completely random direction \tilde{v}_k is just as good as querying according to what the adaptive algorithm specifies.

Unfortunately there are further difficulties to overcome with this approach. Namely, suppose that we define each simulated point v_k^{sim} to be the output $U_k \tilde{v}_k$ of a rotation matrix applied to \tilde{v}_k . We would like to take U_k such that $U_k \tilde{v}_k = v_k((\Lambda^i v_j^{\text{sim}})_{i+j < k})$ but this is no longer computable based on $(\Lambda^i \tilde{v}_j)_{i+j < k}$. However, we note that $\Lambda^i v_j^{\text{sim}} = \Lambda^i U_j \tilde{v}_j = U_j \tilde{\Lambda}^i \tilde{v}_j$ where $\tilde{\Lambda} := U_j^\top \Lambda U_j$. This shows that $\Lambda^i v_j^{\text{sim}}$ is computed from the query of \tilde{v}_j , not on the original matrix Λ but on the modified matrix $\tilde{\Lambda}$, together with the matrix U_j . Since we hope that $\tilde{\Lambda}$ has the same law as Λ , then this is good enough for the purposes of simulating the adaptive algorithm. Actually, in order for the induction to work out, it becomes clear that we need to define a sequence of matrices $\Lambda_1, \Lambda_2, \dots, \Lambda_k$, where each Λ_k is related to the previous Λ_{k-1} via $\Lambda_k = U_k^\top \Lambda_{k-1} U_k$, and U_k is chosen such that $v_k^{\text{sim}} = U_k \tilde{v}_k = v_k((\Lambda_{k-1}^i v_j^{\text{sim}})_{i+j < k})$. Then, we must argue that the simulated sequence $v_1^{\text{sim}}, v_2^{\text{sim}}, \dots, v_k^{\text{sim}}$ has the same law as the algorithm's sequence $v_1^{\text{alg}}, v_2^{\text{alg}}, \dots, v_k^{\text{alg}}$.

This last step, however, turns out to be delicate. Indeed, although it is obvious that for a *fixed* orthogonal matrix U' , the law of Λ is the same as the law of $(U')^\top \Lambda U'$, the rotation matrices U_k we choose in the above argument are dependent on the previous queries and oracle outputs, and are hence dependent on Λ itself. In the presence of such dependence, it is not obvious why the law of Λ_k should be the same as the law of Λ , and to address this we prove a conditioning lemma in Section 5.3.2. Once the conditioning lemma is proved, the remainder of the proof follows along the

lines just described, and the details of the induction are carried out in Section 5.3.3.

3 A general sampling lower bound in dimension two

3.1 Overview

Our goal is to show the following theorem:

Theorem 4 (lower bound in dimension two). *There is a universal constant $\varepsilon_0 > 0$ such that the following holds. The query complexity of sampling from the class of distributions $\pi \propto \exp(-V)$ on \mathbb{R}^2 such that V is 1-strongly convex, κ -smooth, and minimized at 0, with accuracy ε_0 in total variation distance, is at least $\Omega(\log \kappa)$.*

The strategy to do so will be to construct a finite family \mathcal{S} of potentials in the given class which satisfies the following two properties:

- The potentials are *hard to identify via queries* (in the sense of Definition 9 below), and therefore any algorithm must query V at $\Omega(\log \kappa)$ points in order to identify which $V \in \mathcal{S}$ the algorithm is querying.
- The potentials are *well-separated* (in the sense of Definition 10 below), which loosely means that they have mostly non-overlapping support and hence (by Proposition 11) a single sample from $\pi \propto \exp(-V)$ suffices to identify $V \in \mathcal{S}$ with constant probability.

Before describing the potentials \mathcal{S} in more detail, we note some basic definitions.

Definition 5. *Given two functions $f, g : \mathbb{R}^d \rightarrow \mathbb{R}$, the convolution $f * g$ is the function defined as $(f * g)(x) := \int_{\mathbb{R}^d} f(y) g(x - y) dy$, for all $x \in \mathbb{R}^d$.*

Definition 6. *For $\delta > 0$, we define χ_δ to be the indicator function of the ball B_δ of radius δ around the origin. By this, we mean $\chi_\delta(x) = 1$ if $\|x\|_2 \leq \delta$, and $\chi_\delta(x) = 0$ otherwise.*

The family \mathcal{S} of potentials will have cardinality $\kappa^{\Omega(1)}$, so that identification of the potential requires $\Omega(\log \kappa)$ bits of information. Actually, by rescaling the potentials, it suffices for each potential V to be $\kappa^{-O(1)}$ -convex and $\kappa^{O(1)}$ -smooth. Our eventual construction also satisfies the following properties.

- Each $V \in \mathcal{S}$ is of the form $V = \tilde{V} * \chi_\delta + \|\cdot\|^2 / (2\kappa^{O(1)})$, where $\tilde{V} : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a convex, non-negative, and piecewise linear potential, and δ will have scale $\delta = \kappa^{-\Theta(1)}$.
- Each $V \in \mathcal{S}$ is zero in a small neighbourhood of a ray ℓ emanating from the origin, and grows fast outside of this ray; hence, the potentials are well-separated.
- Suppose that ℓ, ℓ' are the rays corresponding to two potentials $V, V' \in \mathcal{S}$. At distances from ℓ and ℓ' that are much larger than the angle $\angle(\ell, \ell')$, the potentials V, V' are exactly equal. This is the property makes the potentials hard to identify via queries.

Throughout the proof, we assume that κ is sufficiently large, $\kappa \geq \Omega(1)$.

3.2 Definitions and the information-theoretic argument

Definition 7 (density and normalizing constant). *Given a strictly convex function $V : \mathbb{R}^d \rightarrow \mathbb{R}$, we denote by P_V the probability distribution with density $Z^{-1} \exp(-V)$ w.r.t. Lebesgue measure, where $Z := \int \exp(-V)$ is the normalizing constant. In an abuse of notation, we also use P_V to refer to the density itself.*

Definition 8 (queries and extended oracle). *For a fixed potential V , and given a query $x \in \mathbb{R}^d$, the extended oracle responds with $V(B_\delta(x))$, which consists of the value of V for all points in the ball of radius δ centered at x . For a sequence of (possibly adaptive and randomized) queries x_1, \dots, x_n and observations $V(B_\delta(x_1)), \dots, V(B_\delta(x_n))$, we denote the information from the i -th query by $\xi_i := \{x_i, V(B_\delta(x_i))\}$, and the information from all the queries by*

$$\xi_{1:n} := \{\xi_1, \dots, \xi_n\}.$$

Note that the extended oracle in Definition 8 provides more information (the set of values of the potential in some ball around the query point x) to the algorithm than our original first-order query model, from which the algorithm only observes $(V(x), \nabla V(x))$ at the query x . A lower bound for sampling in this stronger query model clearly implies a lower bound in the original query model. We consider the stronger model out of technical convenience, as this notion is robust to the mollification in the construction of the potentials.

Definition 9 (hard to identify via queries). *A finite set \mathcal{S} of potentials in \mathbb{R}^d is called \mathcal{I} -hard to identify with queries at scale δ if the following holds: for $V \sim \text{uniform}(\mathcal{S})$, any sequence of queries x_1, \dots, x_n to the extended oracle made by a deterministic adaptive algorithm satisfies*

$$I(\xi_{1:n}; V) \leq \mathcal{I}n,$$

where I denotes the mutual information.

Definition 10 (well-separated set). *A set \mathcal{S} of potentials is well-separated if there is a family of measurable sets $(\Omega_V)_{V \in \mathcal{S}}$ where the sets Ω_V are disjoint, and a universal constant $c > 0$ such that*

$$P_V(\Omega_V) \geq c, \quad \text{for all } V \in \mathcal{S}.$$

The motivation for this definition is the following fact:

Proposition 11 (one sample identifies well-separated distributions). *Let \mathcal{S} be a well-separated set of potentials and conditionally on $V \sim \text{uniform}(\mathcal{S})$, suppose that X is a sample from a probability measure \hat{P}_V which is at most $\frac{c}{2}$ away from P_V in total variation distance. Then,*

$$\mathbb{P}\{X \in \Omega_V\} \geq \frac{c}{2}.$$

Proof. By conditioning on V ,

$$\mathbb{P}\{X \in \Omega_V\} = \mathbb{E} \mathbb{P}\{X \in \Omega_V \mid V\} = \mathbb{E} \hat{P}_V(\Omega_V) \geq \mathbb{E}[P_V(\Omega_V) - \|P_V - \hat{P}_V\|_{\text{TV}}] \geq \frac{c}{2},$$

which is what we wanted to show. \square

This shows that the minimum-distance estimator

$$\hat{V} := \arg \min_{V \in \mathcal{S}} \inf_{z \in \Omega_V} \|X - z\| \tag{3.1}$$

succeeds at estimating the randomly drawn V with constant probability. On the other hand, we have Fano's inequality from information theory.

Theorem 12 (Fano's inequality, [CT06, Theorem 2.10.1]). *Suppose that \mathcal{S} is a finite set and $V \sim \text{uniform}(\mathcal{S})$. Suppose that \hat{V} is any estimator which is based on some data ξ . Then,*

$$\mathbb{P}\{\hat{V} \neq V\} \geq 1 - \frac{I(\xi; V) + \log 2}{\log |\mathcal{S}|}.$$

Fano's inequality enables us to reduce Theorem 4 to the following proposition:

Proposition 13 (well-separated set which is hard to identify via queries). *Let $\kappa \geq \Omega(1)$. Then, there is a set \mathcal{S} of potentials such that:*

1. *All elements of \mathcal{S} are $\kappa^{-O(1)}$ -convex and $\kappa^{O(1)}$ -smooth, and have their minimum at zero.*
2. *\mathcal{S} has cardinality $\kappa^{\Omega(1)}$.*
3. *\mathcal{S} is well-separated with $c = \Omega(1)$.*
4. *\mathcal{S} is hard to identify via queries at scale $\delta = \kappa^{-\Theta(1)}$, and with $\mathcal{I} = O(1)$.*

Proof. [Proof of Theorem 4] Suppose that there is a sampling algorithm which, given any target distribution $\pi \propto \exp(-V)$ on \mathbb{R}^2 such that V is 1-strongly convex, $\bar{\kappa}$ -smooth, and minimized at 0, outputs a sample X whose law is ε_0 close in total variation distance to π using $n(\bar{\kappa})$ queries to the extended oracle. Let \mathcal{S} be the family in Proposition 13. By choosing $\varepsilon_0 = c/2 = \Omega(1)$ and rescaling the potentials accordingly, then Proposition 11 implies that the sampling algorithm can identify $V \sim \text{uniform}(\mathcal{S})$ using $n(\bar{\kappa})$ queries with constant probability, where $\bar{\kappa} = \kappa^{O(1)}$. Namely, for the estimator \hat{V} in (3.1),

$$\mathbb{P}\{\hat{V} = V\} \geq \frac{c}{2} = \Omega(1). \quad (3.2)$$

On the other hand, we can prove a lower bound for the error probability of any estimator \hat{V} constructed using adaptive queries. First we assume that the estimator is deterministic given previous queries. Because the set \mathcal{S} is hard to identify, by Fano's inequality (Theorem 12) we have

$$\mathbb{P}\{\hat{V} \neq V\} \geq 1 - \frac{I(\xi_{1:n(\bar{\kappa})}; V) + \log 2}{\log |\mathcal{S}|} \geq 1 - \frac{\mathcal{I}n(\bar{\kappa}) + \log 2}{\log |\mathcal{S}|} = 1 - O\left(\frac{n(\bar{\kappa})}{\log \kappa}\right), \quad (3.3)$$

for all $n(\bar{\kappa}) \leq c|\mathcal{S}| = O(\log \kappa)$. If the estimator is instead randomized, it depends on a random seed ζ that is independent of V . In this case, the same argument as above conditional on ζ gives

$$\mathbb{P}\{\hat{V} \neq V \mid \zeta\} \geq 1 - \Omega\left(\frac{n(\bar{\kappa})}{\log \kappa}\right).$$

Taking expectation over ζ , we see that (3.3) holds also for randomized algorithms. Combined with (3.2), we see that $n(\bar{\kappa}) \geq \Omega(\log \kappa) = \Omega(\log \bar{\kappa})$. \square

3.3 Reductions and properties of the construction

Recall from Section 3.1 that each $V \in \mathcal{S}$ is of the form $V = \tilde{V} * \chi_\delta + \|\cdot\|^2/(2\kappa^{O(1)})$. In this section, we reduce the desired properties of \mathcal{S} , namely that \mathcal{S} is well-separated and hard to identify via queries, to geometric properties of the potentials summarized in Proposition 14 below.

By increasing κ by a factor of at most two, which will not harm the final lower bound, we can assume that $\kappa = 2^N$ for some positive integer N . We also set $\delta := \kappa^{-5}$. Let B_N denote the set of binary strengths of length N . For each $b \in B_N$ and $\ell \in [N]$, we let $[b]_\ell := 0.00b_1 \dots b_\ell$ in binary representation, and set $[b] := [b]_N$.

Proposition 14 (geometric properties). *There are functions \tilde{V}_b , for $b \in B_N$, such that:*

(P0) \tilde{V}_b is convex and $\kappa^{O(1)}$ -smooth on average at scale $\delta = \kappa^{-5}$, i.e., $\tilde{V}_b * \chi_\delta$ is $\kappa^{O(1)}$ -smooth, and attains its minimum $V_b(0) = 0$ at zero.

(P1) The zero set $\mathcal{Z}_b := \{\tilde{V}_b = 0\}$ contains the $10^3\delta$ -neighborhood of the set

$$\tilde{\mathcal{Z}}_b := \{(x, \beta x) \in \mathbb{R}^2 \mid x \geq 0, [b] - 2^{-N} \leq \beta \leq [b] + 2^{-N}\}, \quad (3.4)$$

and is contained in the 1-neighbourhood of $\tilde{\mathcal{Z}}_b$.

(P2) Moreover, for all $x, y \in \mathbb{R}^2$,

$$\tilde{V}_b(x, y) \geq \kappa^4 (\text{dist}((x, y), \tilde{\mathcal{Z}}_b) - 1)_+.$$

(P3) If b, b' coincide in the first ℓ bits then \tilde{V}_b and $\tilde{V}_{b'}$ coincide in the set

$$\{(x, y) \in \mathbb{R}^2 \mid x < \frac{1}{4} 2^{-3N} \text{ or } |y - [b]_\ell x| > 100 \cdot 2^{-\ell} x\}.$$

We check that these properties imply that Proposition 13 holds.

Proof. [Proof of Proposition 13] Let \mathcal{S} be the collection of potentials $V_b := \tilde{V}_b * \chi_\delta + \|\cdot\|^2/(2\kappa^{16})$ for $b \in B_N$, where $\{\tilde{V}_b : b \in B_N\}$ are the functions from Proposition 14. We now verify the four properties of Proposition 13.

Proof of 1. By (P0), we know that \tilde{V}_b is convex, which implies that $\tilde{V}_b * \chi_\delta$ is also convex. Therefore, V_b is κ^{-16} -strongly convex. In addition, by (P0), $\tilde{V}_b * \chi_\delta$ is $\kappa^{O(1)}$ -smooth, which means that V_b is $\kappa^{O(1)} + \kappa^{-16} \leq \kappa^{O(1)}$ -smooth.

Proof of 2. By construction, $|\mathcal{S}| = \kappa$.

Proof of 3. We now show that \mathcal{S} is c -separated. For any string b , recall the definition of $\tilde{\mathcal{Z}}_b$ from (3.4). Define the set

$$\Omega_b := \{(x, \beta x) \in \mathbb{R}^2 \mid x \geq 2^{-3N}, [b] - 0.4 \cdot 2^{-N} \leq \beta \leq [b] + 0.4 \cdot 2^{-N}\}.$$

It is clear that $\{\Omega_b : b \in B_N\}$ is a family of disjoint sets. By (P1) we know that the zero set \mathcal{Z}_b of \tilde{V}_b contains a $10^3\delta$ -neighborhood of $\tilde{\mathcal{Z}}_b$. Since $\Omega_b \subset \tilde{\mathcal{Z}}_b$, it follows that $\tilde{V}_b * \chi_\delta = 0$ on Ω_b .

Let $\tilde{\Omega}_b := \{(x, y) \in \Omega_b : \|(x, y)\| \leq \kappa^8\}$. Note that the full set of points (x, y) with $\|(x, y)\| \leq \kappa^8$ has volume $\pi\kappa^{16}$, and Ω_b is a sector of the plane with arc $\Theta(2^{-N})$, minus a small set of points (specifically, the points in the sector with $x \leq 2^{-3N}$, which also means $y \leq O(2^{-3N})$). Therefore, the volume of $\tilde{\Omega}_b$ is $\Theta(\kappa^{16} \cdot 2^{-N}) = \Theta(\kappa^{15})$. In addition, all points $(x, y) \in \tilde{\Omega}_b$ have $V_b(x, y) = -\|(x, y)\|^2/(2\kappa^{16}) \geq -1/2$. Hence,

$$\int_{\Omega_b} \exp(-V_b) \geq \int_{\tilde{\Omega}_b} \exp(-V_b) \geq \Omega(\kappa^{15}). \quad (3.5)$$

Next, we bound the full integral of $\exp(-V_b)$ across \mathbb{R}^d by splitting \mathbb{R}^d into four regions $\mathbb{R}^d = \tilde{\mathcal{Z}}_b \cup \Psi_{1,b} \cup \Psi_{2,b} \cup \Psi_{3,b}$, defined as follows:

- $\Psi_{1,b} := \{(x, y) \in \mathbb{R}^2 \setminus \tilde{\mathcal{Z}}_b : \text{dist}((x, y), \tilde{\mathcal{Z}}_b) \leq 2, \|(x, y)\| \leq \kappa^9\}.$
- $\Psi_{2,b} := \{(x, y) \in \mathbb{R}^2 \setminus (\tilde{\mathcal{Z}}_b \cup \Psi_{1,b}) : \|(x, y)\| \leq \kappa^9\}.$

- $\Psi_{3,b} = \mathbb{R}^2 \setminus (\tilde{\mathcal{Z}}_b \cup \Psi_{1,b} \cup \Psi_{2,b})$.

Note that all points $\Psi_{3,b}$ have norm at least κ^9 . To show that most of the mass of P_{V_b} is concentrated on $\tilde{\mathcal{Z}}_b$, we must show that the integrals over $\Psi_{1,b}$, $\Psi_{2,b}$, and $\Psi_{3,b}$ are small. In a nutshell, the integral over $\Psi_{1,b}$ is small because the 2-neighborhood of $\tilde{\mathcal{Z}}_b$ is small (relative to the size of $\tilde{\mathcal{Z}}_b$ itself); the integral over $\Psi_{2,b}$ is small because \tilde{V}_b increases rapidly outside $\tilde{\mathcal{Z}}_b$; and the integral over $\Psi_{3,b}$ is small because the Gaussian part of V_b is small over this region.

On these four regions, we have the following bounds. First, $\int_{\mathbb{R}^2} \exp(-\|\cdot\|^2/(2\kappa^{16})) = 2\pi\kappa^{16}$. Therefore, since the sector $\tilde{\mathcal{Z}}_b$ has arc $\Theta(2^{-N})$, by rotational symmetry

$$\int_{\tilde{\mathcal{Z}}_b} \exp(-V_b) \leq \int_{\tilde{\mathcal{Z}}_b} \exp(-\frac{\|\cdot\|^2}{2\kappa^{16}}) \leq O(2^{-N}) \int_{\mathbb{R}^2} \exp(-\frac{\|\cdot\|^2}{2\kappa^{16}}) \leq O(\kappa^{15}).$$

Note that $\Psi_{1,b}$ consists of two strips adjacent to $\tilde{\mathcal{Z}}_b$, where each strip has width 2 and length $O(\kappa^9)$, together with a piece of area $O(1)$ near the origin. Thus, $\text{vol}(\Psi_{1,b}) \leq O(\kappa^9)$, yielding

$$\int_{\Psi_{1,b}} \exp(-V_b) \leq \text{vol}(\Psi_{1,b}) \leq O(\kappa^9).$$

Next, for $(x, y) \in \mathbb{R}^2$ such that $\text{dist}((x, y), \tilde{\mathcal{Z}}_b) \geq 3/2$, by (P2) we have $\tilde{V}_b(x, y) \geq \kappa^4$. After mollification at scale $\delta \leq 1/2$, we conclude that $\tilde{V}_b * \chi_\delta \geq \kappa^4$ on $\Psi_{2,b}$. In addition, $\Psi_{2,b}$ is contained in the ball of radius κ^9 , so the volume of $\Psi_{2,b}$ is at most $\pi\kappa^{18}$. Therefore,

$$\int_{\Psi_{2,b}} \exp(-V_b) \leq \pi\kappa^{18} \exp(-\kappa^4).$$

Finally, all points in $\Psi_{3,b}$ have ℓ_2 norm at least κ^9 , so

$$\int_{\Psi_{3,b}} \exp(-V_b) \leq \iint_{\|\cdot\| \geq \kappa^9} \exp(-\frac{\|\cdot\|^2}{2\kappa^{16}}) \leq O(\kappa^8) \exp(-\Omega(\kappa^2)),$$

by standard Gaussian tail estimates. Therefore,

$$\int_{\mathbb{R}^2} \exp(-V_b) \leq O(\kappa^{15} + \kappa^9 + \exp(-\Omega(\kappa^4)) + \exp(-\Omega(\kappa^2))) \leq O(\kappa^{15}). \quad (3.6)$$

Overall, (3.5) and (3.6) together imply that $P_{V_b}(\Omega_b) \geq \Omega(1)$, i.e., \mathcal{S} is $\Omega(1)$ -well-separated.

Proof of 4. Finally, we show that \mathcal{S} is hard to identify via queries at scale $\delta = \kappa^{-\Theta(1)}$ with $\mathcal{I} = O(1)$. We consider b drawn uniformly at random from B_N .

First, however, we need to extend (P3) to V_b (i.e., taking into account the mollification at scale δ). We claim that if b, b' coincide in the first ℓ bits, then V_b and $V_{b'}$ coincide in the set

$$\{(x, y) \in \mathbb{R}^2 \mid x < \frac{1}{8} 2^{-3N} \text{ or } |y - [b]_\ell x| > 200 \cdot 2^{-\ell} x\}. \quad (3.7)$$

In light of (P3), it suffices to show that if (x, y) lies in this set and $\|(x', y') - (x, y)\| \leq \delta$, then $x' < \frac{1}{4} 2^{-3N}$ or $|y' - [b]_\ell x'| > 100 \cdot 2^{-\ell} x'$. In other words, the δ -neighborhood of (3.7) is contained in the set in (P3). In the first case, $x' < \frac{1}{4} 2^{-3N}$ follows if $\delta < \frac{1}{8} 2^{-3N}$, but since $\delta = \kappa^{-5} = 2^{-5N}$ this holds for large κ . In the second case,

$$|y' - [b]_\ell x'| \geq |y - [b]_\ell x| - \delta - [b]_\ell \delta \geq 200 \cdot 2^{-\ell} x - 2\delta.$$

This is greater than $100 \cdot 2^{-\ell} x$ provided that $2\delta \leq 100 \cdot 2^{-\ell} x$, but this follows because $\delta = 2^{-5N}$ and $x \geq \frac{1}{8} 2^{-3N}$ (as we are in the negation of the first case). In fact, by replacing δ with 2δ , the same argument shows that for all (x, y) lying in the set (3.7), we have $V_b(B_\delta(x, y)) = V_{b'}(B_\delta(x, y))$. Note also that (3.7) shows that it is useless to query any points (x, y) with $x < \frac{1}{8} 2^{-3N}$, so for the remainder of the proof we assume that the algorithm does not do so.

We now move to a stronger oracle model. Namely, given a query point $(x, y) \in \mathbb{R}^2$, let ℓ be the largest integer such that $|y - [b]_\ell x| \leq 200 \cdot 2^{-\ell} x$. Then, the oracle outputs $\hat{\xi} := [b]_{\ell+1}$, i.e., the oracle reveals the first $\ell + 1$ bits of b . To see that this new oracle is indeed stronger, observe that we can simulate the previous oracle using the revealed bits $[b]_{\ell+1}$; namely, pick any bit string b' which is consistent, in the sense that $[b']_{\ell+1} = [b]_{\ell+1}$. Then, by the choice of ℓ , we have $|y - [b]_{\ell+1} x| > 200 \cdot 2^{-(\ell+1)} x$, so that $V_b(B_\delta(x, y)) = V_{b'}(B_\delta(x, y))$, and hence we can output $V_b(B_\delta(x, y))$ given knowledge of $[b]_{\ell+1}$. It therefore suffices to bound the mutual information $I(\hat{\xi}_{1:n}; b)$ where $\hat{\xi}_{1:n}$ denotes the output of the stronger oracle on a sequence of adaptive but deterministic queries $(x_1, y_1), \dots, (x_n, y_n)$.

We can then write

$$\begin{aligned} I(\hat{\xi}_{1:n}; b) &= \sum_{i=1}^n I(\hat{\xi}_i; b \mid \hat{\xi}_{1:i-1}) \\ &= \sum_{i=1}^n \{H(\hat{\xi}_i \mid \hat{\xi}_{1:i-1}) - H(\hat{\xi}_i \mid \hat{\xi}_{1:i-1}, b)\} \\ &\leq \sum_{i=1}^n H(\hat{\xi}_i \mid \hat{\xi}_{1:i-1}), \end{aligned} \tag{3.8}$$

where $H(\cdot \mid \cdot)$ denotes the conditional entropy. The first line follows from the chain rule for mutual information, the second line follows from definition of mutual information, and third line follows from non-negativity of conditional entropy. Thus, we are done if we can show that $H(\hat{\xi}_i \mid \hat{\xi}_{1:i-1}) \leq O(1)$, for all $i \leq c|\mathcal{S}|$.

Conditionally on any particular realization of $\hat{\xi}_{1:i-1}$, let ℓ_0 denote the number of bits of b revealed thus far and let $[b_0]_{\ell_0}$ denote the revealed bits. Clearly the bit string b is uniformly distributed on the set B'_N of bit strings b' with $[b']_{\ell_0} = [b_0]_{\ell_0}$. Also, since we have assumed that the algorithm's queries are deterministic given the past history, the next query point (x_i, y_i) is deterministic. Then, the conditional probability that $\ell \geq \ell_0$ bits are revealed by the next query is

$$\begin{aligned} &\mathbb{P}\{200 \cdot 2^{-\ell} x_i < |y_i - [b]_\ell x_i| \leq 200 \cdot 2^{-(\ell-1)} x_i \mid \hat{\xi}_{1:i-1}\} \\ &\leq \mathbb{P}\left\{\frac{y_i}{x_i} - 200 \cdot 2^{-(\ell-1)} \leq [b]_\ell \leq \frac{y_i}{x_i} + 200 \cdot 2^{-(\ell-1)} \mid \hat{\xi}_{1:i-1}\right\}. \end{aligned}$$

This is the probability that a uniformly chosen element of B'_N belongs to an interval of length $\Theta(2^{-\ell})$. Since there are $2^{N-\ell_0}$ elements of B'_N , and $\Theta(2^{N-\ell})$ of them belong to any fixed interval of length $\Theta(2^{-\ell})$, we conclude that the above probability is $O(2^{-(\ell-\ell_0)})$.

We then have

$$H(\hat{\xi}_i \mid \hat{\xi}_{1:i-1}) \leq \mathbb{E} \sum_{\ell \geq \ell_0} (\ell - \ell_0) O(2^{-(\ell-\ell_0)}) \leq O(1),$$

where the expectation is taken over ℓ_0 (which depends on the realization of $\hat{\xi}_{1:i-1}$). Substituting the above bound into (3.8), we conclude that $I(\xi_{1:n}; b) = O(n)$, which implies that \mathcal{S} is indeed hard to identify via queries. \square

3.4 Construction of the distributions

This section contains the proof of Proposition 14.

For integers $1 \leq k \leq N$, let $[b]_k$ be the number $0.00b_1b_2 \dots b_k$ in binary representation, and let $[b]_k := [b] := [b]_N$ for $k \geq N$. Define

$$\phi_{k,b}(x, y) := (|y - [b]_k x| - (2^{-k}x + 2^{-(3N-k)}))_+. \quad (3.9)$$

Here, the term $2^{-(3N-k)}$ essentially controls the thickness of the slab, and in particular, the slab becomes thicker for larger k ; this ensures that the maximum of the $\phi_{k,b}$ will be dominated by small k far away. We also write $\phi_k := \phi_{k,b}$ when b is clear from context. For $x \geq 0$, the function ϕ_k essentially measures the distance to the set

$$\{(x, [b]_k x + \xi_k) \in \mathbb{R}^2 : x \geq 0, |\xi_k| \leq 2^{-k}x + 2^{-(3N-k)}\}.$$

Finally, we define the potential

$$\tilde{V}_b(x, y) := 2^{7N} \max_{k=1, \dots, N} 2^{-k} \phi_k(x, y). \quad (3.10)$$

Proof. [Proof of Proposition 14] We prove that the construction (3.10) satisfies each of the four properties in turn.

Proof of Property (P0). The convexity of \tilde{V}_b follows because each ϕ_k is convex. To check that \tilde{V}_b is $\kappa^{O(1)}$ -smooth on average, using the compositionality of the maximum (i.e., $\max(a, \max(b, c)) = \max(a, b, c)$) we see that that \tilde{V}_b can be written as a maximum of affine functions, each of slope $\kappa^{O(1)}$; hence, \tilde{V}_b is $\kappa^{O(1)}$ -Lipschitz. Differentiating under the integral,

$$\nabla(\tilde{V}_b * \chi_\delta)(x, y) = \iint_{B_\delta} \nabla \tilde{V}_b(x + u, y + v) du dv = \iint \nabla \tilde{V}_b \mathbb{1}_{B_\delta(x, y)},$$

where the expression makes sense because \tilde{V}_b is Lipschitz and hence differentiable a.e. by Rademacher's theorem, and the absolute continuity of \tilde{V}_b ensures the validity of the fundamental theorem of calculus. Then, by Hölder's inequality,

$$\begin{aligned} \|\nabla(\tilde{V}_b * \chi_\delta)(x, y) - \nabla(\tilde{V}_b * \chi_\delta)(x', y')\| &\leq (\sup \|\nabla \tilde{V}_b\|) \|\mathbb{1}_{B_\delta(x, y)} - \mathbb{1}_{B_\delta(x', y')}\|_{L^1} \\ &\leq \kappa^{O(1)} \text{vol}(B_\delta(x, y) \triangle B_\delta(x', y')). \end{aligned}$$

By elementary considerations, the volume of the symmetric difference between the balls is bounded by $O(\kappa^{O(1)} \|(x, y) - (x', y')\|)$, and therefore $\nabla(\tilde{V}_b * \chi_\delta)$ is $\kappa^{O(1)}$ -Lipschitz.

Finally, it is obvious that $\tilde{V}_b \geq 0$ and $\tilde{V}_b = 0$ at the origin.

Proof of Property (P1). We only need to verify that any point (x, y) which is $10^3\delta$ -close to \tilde{Z}_b satisfies $\tilde{V}_b(x, y) = 0$, as the second part of Property (P1) is automatically implied by Property (P2). For such a point (x, y) , there exists (x', y') such that

$$x' \geq 0, \quad |x' - x| \wedge |y' - y| \leq 10^3\delta, \quad \text{and} \quad |y' - [b]x'| \leq 2^{-N}x'.$$

This also implies $|y' - [b]_k x'| \leq 2^{-k}x'$ for all $1 \leq k \leq N$, since $|[b]_k - [b]| \leq 2^{-k} - 2^{-N}$. Therefore, for all $1 \leq k \leq N$, $|y - [b]_k x| \leq 2^{-k}(x + 10^3\delta) + 2 \cdot 10^3\delta \leq 2^{-k}x + 2^{-(3N-k)}$, since $\delta = 2^{-5N}$. By the definition (3.10) of \tilde{V}_b and the definition of ϕ_k in (3.9), it follows that $\tilde{V}_b(x, y) = 0$.

Proof of Property (P2). We just need to check that

$$2^{6N} \phi_N(x, y) \geq \kappa^4 (\text{dist}((x, y), \tilde{\mathcal{Z}}_b) - 1)_+,$$

or equivalently, $2^{2N} \phi_N(x, y) \geq (\text{dist}((x, y), \tilde{\mathcal{Z}}_b) - 1)_+$. We first consider the case when $x \geq 0$, and we may assume that $(x, y) \notin \tilde{\mathcal{Z}}_b$ as otherwise the claim is obvious. If (x, y) has distance Δ to its closest point in $\tilde{\mathcal{Z}}_b$, then any y' such that $(x, y') \in \tilde{\mathcal{Z}}_b$ must satisfy $|y - y'| \geq \Delta$. Applying this to $y' = [b]x \pm 2^{-N}x$, we obtain

$$\text{dist}((x, y), \tilde{\mathcal{Z}}_b) \leq |y - [b]x + 2^{-N}x| \wedge |y - [b]x - 2^{-N}x| = |y - [b]x| - 2^{-N}x.$$

In turn, it implies that $\phi_N(x, y) \geq (\text{dist}((x, y), \tilde{\mathcal{Z}}_b) - 2^{-(3N-k)})_+ \geq (\text{dist}((x, y), \tilde{\mathcal{Z}}_b) - 1)_+$.

If $x < 0$, then $\text{dist}((x, y), \tilde{\mathcal{Z}}_b) \leq \|(x, y)\| \leq \sqrt{2} \max(|x|, |y|)$. Then, for N large,

$$\begin{aligned} 2^{2N} \phi_N(x, y) &= 2^{2N} (|y - [b]x| - 2^{-N}x - 2^{-(3N-k)})_+ \\ &= 2^{N-1/2} (2^{N+1/2} |y - [b]x| + \sqrt{2}|x| - 2^{-(2N-k)+1/2})_+ \\ &\geq 2^{N-1/2} \left(2^{3/2} \max(0, |y| - \frac{1}{2}|x|) + \sqrt{2}|x| - 1 \right)_+ \\ &\geq 2^{N-1/2} (\sqrt{2} \max(|x|, |y|) - 1)_+ \geq (\text{dist}((x, y), \tilde{\mathcal{Z}}_b) - 1)_+. \end{aligned}$$

The first inequality follows because $|y - [b]x| = ||y| - [b] \text{sgn}(y)x| \geq |y| - \frac{1}{2}|x| + \frac{1}{2}|x| - [b] \text{sgn}(y)x \geq |y| - \frac{1}{2}|x|$ and because N is sufficiently large.

Proof of Property (P3). The last property follows from Proposition 15 below, because if b, b' agree on the first ℓ bits, then on the set in the statement of Property (P3),

$$\tilde{V}_b = 2^{7N} \max_{k=1, \dots, N} 2^{-k} \phi_{k,b} = 2^{7N} \max_{k=1, \dots, \ell} 2^{-k} \phi_{k,b} = 2^{7N} \max_{k=1, \dots, \ell} 2^{-k} \phi_{k,b'} = 2^{7N} \max_{k=1, \dots, N} 2^{-k} \phi_{k,b'} = \tilde{V}_{b'}.$$

The second and fourth equalities invoke Proposition 15, and the third equality uses the fact that $\phi_{k,b}$ only depends on b through $[b]_k$. This completes the proof. \square

Proposition 15 (potentials agree if bits agree). *Let $S_\ell(b)$ be the set*

$$S_\ell(b) := \{(x, y) \in \mathbb{R}^2 : x < \frac{1}{4} 2^{-3N} \text{ or } |y - [b]_\ell x| \geq 100 \cdot 2^{-\ell} x\}.$$

Then, for $x, y \in S_\ell(b)$,

$$\max_{k=1, \dots, N} 2^{-k} \phi_k(x, y) = \max_{k=1, \dots, \ell} 2^{-k} \phi_k(x, y).$$

In turn, Proposition 15 follows by induction from:

Proposition 16 (induction). *If $(x, y) \in S_\ell(b)$, and for some $k > \ell$ we have $\phi_k(x, y) > 0$, then $\phi_k(x, y) \leq 2\phi_{k-1}(x, y)$.*

Proof. First, we may assume that $x > 0$. This is because if $x \leq 0$,

$$\begin{aligned} \phi_{k-1}(x, y) &\geq |y - [b]_k x| - |[b]_{k-1} - [b]_k| |x| - 2^{-(k-1)}x - 2^{-(3N-k+1)} \\ &\geq |y - [b]_k x| + 2^{-k}x - 2^{-(k-1)}x - 2^{-(3N-k+1)} \\ &= |y - [b]_k x| - 2^{-k}x - 2^{-(3N-k+1)} \\ &\geq \phi_k(x, y), \end{aligned}$$

since we are assuming $\phi_k(x, y) > 0$.

Now, since $x > 0$, we start by estimating

$$\begin{aligned}\phi_{k-1}(x, y) &\geq |y - [b]_k x| - |[b]_{k-1} - [b]_k| x - 2^{-(k-1)}x - 2^{-(3N-k+1)} \\ &\geq |y - [b]_k x| - 3 \cdot 2^{-k}x - 2^{-(3N-k+1)} \\ &= \phi_k(x, y) - 2 \cdot 2^{-k}x + 2^{-(3N-k+1)}\end{aligned}$$

and

$$\phi_k(x, y) = |y - [b]_k x| - (2^{-k}x + 2^{-(3N-k)}).$$

First, suppose that $x \leq \frac{1}{4} 2^{-3N}$. Then, $2^{-(3N-k+1)} \geq 2 \cdot 2^{-k}x$, so in fact $\phi_{k-1}(x, y) \geq \phi_k(x, y)$. Alternatively, if $x \geq \frac{1}{4} 2^{-3N}$ and $|y - [b]_\ell x| \geq 100 \cdot 2^{-\ell}x$, then

$$\begin{aligned}2\phi_{k-1}(x, y) &\geq 2|y - [b]_\ell x| - 2|[b]_\ell - [b]_{k-1}|x - 4 \cdot 2^{-k}x - 2^{-(3N-k)} \\ &\geq 2|y - [b]_\ell x| - 6 \cdot 2^{-\ell}x - 2^{-(3N-k)}, \\ \phi_k(x, y) &\leq |y - [b]_\ell x| + |[b]_\ell - [b]_k|x - 2^{-k}x - 2^{-(3N-k)} \\ &\leq |y - [b]_\ell x| + 2^{-\ell}x - 2^{-(3N-k)}.\end{aligned}$$

As a result, when $|y - [b]_\ell x| \geq 100 \cdot 2^{-\ell}x$, we see that $\phi_{k-1}(x, y) \geq \frac{1}{2} \phi_k(x, y)$. \square

4 A lower bound for sampling from Gaussians via Wishart matrices

We define $W \sim \text{Wishart}(d)$ to mean $W = XX^\top$ where each entry of $X \in \mathbb{R}^{d \times d}$ is $\mathcal{N}(0, \frac{1}{d})$. We aim to prove the following two theorems, which together imply a query complexity lower bound for sampling from Gaussians.

Theorem 17 (reducing inverse trace estimation to sampling). *Let $\delta > 0$. There is a universal constant $c > 0$ (depending only on δ) such that the following hold. Suppose that $d \geq c^{-1}$ and there exists a query algorithm such that, for any Gaussian target distribution $\pi := \mathcal{N}(0, \Sigma)$ in \mathbb{R}^d with $cd^{-2} I_d \preceq \Sigma^{-1} \preceq c^{-1} I_d$, outputs a sample from a distribution $\hat{\pi}$ such that either $\|\hat{\pi} - \pi\|_{\text{TV}} \leq c$ or $\sqrt{cd^{-2}} W_2(\hat{\pi}, \pi) \leq c$, using n queries to π .*

Then, given $W \sim \text{Wishart}(d)$, there exists an algorithm which makes at most $c^{-1}n$ matrix-vector queries to W and outputs an estimator $\hat{\text{tr}}$ such that $\frac{1}{2} \text{tr}(W^{-1}) \leq \hat{\text{tr}} \leq 2 \text{tr}(W^{-1})$ with probability at least $1 - \delta$.

Theorem 18 (lower bound for inverse trace estimation). *Let $W \sim \text{Wishart}(d)$ for $d \geq 2$. For any $C > 0$, there exists $\delta > 0$ (depending only on C) such that any algorithm which makes n matrix-vector queries to W and outputs an estimator $\hat{\text{tr}}$ such that $C^{-1} \text{tr}(W^{-1}) \leq \hat{\text{tr}} \leq C \text{tr}(W^{-1})$ with probability at least $1 - \delta$ must use $n \geq \Omega(d)$ queries.*

Remark. Suppose that we want to sample from a target distribution π which is α -strongly log-concave. It is straightforward to check that total variation guarantees are invariant under rescaling the target (replacing π with $S_\# \pi$, where $S : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the scaling map $Sx := \zeta x$ for some $\zeta > 0$), whereas Wasserstein guarantees are not. Instead, the scale-invariant quantity is $\sqrt{\alpha} W_2$, which is what appears in Theorem 17.

Consider the class of centered Gaussian distributions on \mathbb{R}^d which are α -strongly log-concave and β -log-smooth; let $\kappa := \beta/\alpha$ denote the condition number. Let $\mathcal{C}_{\text{G},d}(\kappa, d, \varepsilon)$ denote the query

complexity of outputting a sample which is ε -close in the metric \mathbf{d} to a target distribution in this class, where \mathbf{d} is one of the scale-invariant distances $\mathbf{d} \in \{\text{TV}, \sqrt{\alpha} W_2\}$. Then, Theorems 17 and 18 (with $C = 2$ and δ, c being universal constants) show that for $d \geq c^{-1}$,

$$\mathcal{C}_{\mathbf{G}, \mathbf{d}}(c^{-2}d^2, d, c) \geq \Omega(d). \quad (4.1)$$

By embedding the construction into higher dimensions, we obtain the following corollary.

Corollary 19 (query lower bound via Wishart matrices). *For $\mathbf{d} \in \{\text{TV}, \sqrt{\alpha} W_2\}$, there is a universal constant $c > 0$ such that*

$$\mathcal{C}_{\mathbf{G}, \mathbf{d}}(\kappa, d, c) \geq \Omega(\sqrt{\kappa} \wedge d).$$

Proof. If $\kappa \geq c^{-2}d^2$, then (4.1) yields

$$\mathcal{C}_{\mathbf{G}, \mathbf{d}}(\kappa, d, c) \geq \Omega(d) \geq \Omega(\sqrt{\kappa} \wedge d).$$

Otherwise, if $\kappa \leq c^{-2}d^2$, let d_\star be the largest integer such that $\kappa \geq c^{-2}d_\star^2$. Then, by embedding the d_\star -dimensional construction into dimension d ,

$$\mathcal{C}_{\mathbf{G}, \mathbf{d}}(\kappa, d, c) \geq \mathcal{C}_{\mathbf{G}, \mathbf{d}}(\kappa, d_\star, c) \geq \Omega(d_\star) \geq \Omega(\sqrt{\kappa} \wedge d),$$

which concludes the proof. \square

4.1 Reducing inverse trace estimation to sampling

In this section, we prove Theorem 17, which is based on the concentration of the squared norm of a Gaussian. We recall the following identity:

Lemma 20 (concentration of the squared norm). *Let $Z \sim \mathcal{N}(0, \Sigma)$. Then,*

$$\text{var}(\|Z\|^2) = 2 \|\Sigma\|_{\text{HS}}^2.$$

Proof. Note that since all quantities are rotationally invariant, we may assume without loss of generality that Σ is diagonal. Then the equality claimed is just the variance of a non-homogenous chi-squared random variable. \square

We now prove Theorem 17.

Proof. [Proof of Theorem 17] Let $W \sim \text{Wishart}(d)$ and let $\Sigma := W^{-1}$. By Proposition 23, there exists $c > 0$ (depending only on δ) such that with probability at least $1 - \delta/3$, it holds that

$$cd^{-2} I_d \preceq \Sigma^{-1} \preceq c^{-1} I_d.$$

We work on the event \mathcal{E} that this holds.

Case 1: total variation distance. From Lemma 20 and Chebyshev's inequality, we deduce that if $Z_1, \dots, Z_m \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma)$ and $\widehat{\text{tr}}_\star := m^{-1} \sum_{i=1}^m \|Z_i\|^2$,

$$\mathbb{P}\{|\widehat{\text{tr}}_\star - \text{tr } \Sigma| \geq \frac{1}{2} \text{tr } \Sigma\} \leq \frac{\text{var } \widehat{\text{tr}}_\star}{(\text{tr } \Sigma)^2/4} = \frac{8}{m} \cdot \frac{\text{tr}(\Sigma^2)}{\text{tr}(\Sigma)^2} \leq \frac{8}{m}.$$

Take $m \geq 48/\delta$ so that this probability is at most $\delta/3$. Conditionally on W , let $\hat{\pi}_W$ denote the law of the sample X of the algorithm when run on the target $\mathcal{N}(0, \Sigma)$. By running the sampling algorithm m times, we can obtain i.i.d. samples $X_1, \dots, X_m \stackrel{\text{i.i.d.}}{\sim} \hat{\pi}_W$. Then, for $\hat{\text{tr}} := m^{-1} \sum_{i=1}^m \|X_i\|^2$,

$$\begin{aligned} \mathbb{P}\left\{|\hat{\text{tr}} - \text{tr } \Sigma| \geq \frac{1}{2} \text{tr } \Sigma\right\} &\leq \mathbb{P}(\mathcal{E}^c) + \mathbb{P}\left\{|\hat{\text{tr}} - \text{tr } \Sigma| \geq \frac{1}{2} \text{tr } \Sigma, \mathcal{E}\right\} \\ &\leq \frac{\delta}{3} + \mathbb{E}\left[\mathbb{P}\left\{|\hat{\text{tr}}_\star - \text{tr } \Sigma| \geq \frac{1}{2} \text{tr } \Sigma \mid W\right\} \mathbb{1}_{\mathcal{E}}\right] + \mathbb{E}\left[\|\hat{\pi}_W^{\otimes m} - \mathcal{N}(0, \Sigma)^{\otimes m}\|_{\text{TV}} \mathbb{1}_{\mathcal{E}}\right] \\ &\leq \frac{\delta}{3} + \frac{\delta}{3} + cm. \end{aligned}$$

If we choose $c \leq \delta/(3m)$, then $\hat{\text{tr}}$ is an estimator of $\text{tr}(W^{-1})$ with multiplicative error at most 2 which succeeds with probability at least $1 - \delta$. Note that both c and m depend only on δ .

Case 2: Wasserstein distance. Consider a coupling of X and Z such that, conditionally on W , we have $\mathbb{E}[\|X - Z\|^2 \mid W] = \mathbb{E}[W_2^2(\hat{\pi}_W, \mathcal{N}(0, \Sigma)) \mid W]$. Let $(X_1, Z_1), \dots, (X_m, Z_m)$ be i.i.d. copies of this coupling. Also, let \mathcal{E}' denote the event that $\lambda_{\min}(W^{-1}) \geq \bar{c}d^2$, where \bar{c} is a constant depending only on δ , chosen so that $\mathbb{P}(\mathcal{E}'^c) \leq \delta/3$ using Proposition 23. Then, conditionally on W in the event $\mathcal{E} \cap \mathcal{E}'$,

$$\begin{aligned} \mathbb{E}[|\hat{\text{tr}} - \text{tr } \Sigma| \mid W] &\leq \mathbb{E}[|\hat{\text{tr}} - \hat{\text{tr}}_\star| \mid W] + \mathbb{E}[|\hat{\text{tr}}_\star - \text{tr } \Sigma| \mid W] \\ &\leq \mathbb{E}[|\hat{\text{tr}} - \hat{\text{tr}}_\star| \mid W] + \frac{2 \text{tr } \Sigma}{\sqrt{m}}, \end{aligned}$$

where we used Lemma 20. Using $\|x\|^2 - \|y\|^2 = \langle x - y, x + y \rangle$, for any $\lambda > 0$,

$$\begin{aligned} \mathbb{E}[|\hat{\text{tr}} - \hat{\text{tr}}_\star| \mid W] &\leq \mathbb{E}[|\|X\|^2 - \|Z\|^2| \mid W] \leq \mathbb{E}[\|X - Z\|^2 \mid W] + 2\mathbb{E}[|\langle X - Z, Z \rangle| \mid W] \\ &\leq (1 + \lambda) \mathbb{E}[\|X - Z\|^2 \mid W] + \frac{1}{\lambda} \mathbb{E}[\|Z\|^2 \mid W] \\ &\leq (1 + \lambda) c^3 d^2 + \frac{\text{tr } \Sigma}{\lambda} \leq (1 + \lambda) \frac{c^3}{\bar{c}} \text{tr } \Sigma + \frac{\text{tr } \Sigma}{\lambda}. \end{aligned}$$

For the last line, recall that we are assuming $\mathbb{E}[W_2^2(\hat{\pi}_W, \mathcal{N}(0, \Sigma)) \mid W] \leq c^3 d^2$. If we take $\lambda = 18/\delta$, $m \geq (36/\delta)^2$, and if c is sufficiently small (depending only on δ), we obtain

$$\mathbb{E}[|\hat{\text{tr}} - \text{tr } \Sigma| \mid W] \leq \frac{\delta \text{tr } \Sigma}{6}.$$

By Markov's inequality,

$$\begin{aligned} \mathbb{P}\left\{|\hat{\text{tr}} - \text{tr } \Sigma| \geq \frac{1}{2} \text{tr } \Sigma\right\} &\leq \mathbb{P}(\mathcal{E}^c) + \mathbb{P}(\mathcal{E}'^c) + \mathbb{E}\left[\mathbb{P}\left\{|\hat{\text{tr}} - \text{tr } \Sigma| \geq \frac{1}{2} \text{tr } \Sigma \mid W\right\} \mathbb{1}_{\mathcal{E} \cap \mathcal{E}'}\right] \\ &\leq \frac{\delta}{3} + \frac{\delta}{3} + \frac{\delta}{3} \leq \delta. \end{aligned}$$

We conclude as before. \square

4.2 Lower bound for inverse trace estimation

In this section, we prove Theorem 18. The idea is that due to the heavy tails of $\lambda_{\min}(W^{-1})$ implied by Proposition 23, with some small probability δ , $\text{tr}(W^{-1})$ will be very large. An algorithm for inverse trace estimation which succeeds with probability at least $1 - \delta$ must be able to detect this event, and we show that this requires making $\Omega(d)$ queries.

The key technical tools are the following propositions, due to [BHSW20].

Proposition 21 ([BHSW20, Lemma 3.4]). *Let $W \sim \text{Wishart}(d)$. Then, for any sequence of $n < d$ (possibly adaptive) queries v_1, \dots, v_n and responses $w_1 = Wv_1, \dots, w_n = Wv_n$, there exists an orthogonal matrix $V \in \mathbb{R}^{d \times d}$ and matrices $Y_1 \in \mathbb{R}^{n \times n}, Y_2 \in \mathbb{R}^{(d-n) \times n}$ that only depend on $v_1, \dots, v_n, w_1, \dots, w_n$, such that VWV^\top has the block form*

$$VWV^\top = \begin{bmatrix} Y_1 Y_1^\top & Y_1 Y_2^\top \\ Y_2 Y_1^\top & Y_2 Y_2^\top + \widetilde{W} \end{bmatrix}.$$

Here, conditionally on $v_1, \dots, v_n, w_1, \dots, w_n$, the matrix \widetilde{W} has the $\text{Wishart}(d-n)$ distribution.

Proposition 22 ([BHSW20, Lemma 3.5]). *For any matrices $Y_1 \in \mathbb{R}^{n \times n}, Y_2 \in \mathbb{R}^{(d-n) \times n}$, and any symmetric matrix $\widetilde{W} \in \mathbb{R}^{(d-n) \times (d-n)}$, it holds that*

$$\lambda_{\min} \left(\begin{bmatrix} Y_1 Y_1^\top & Y_1 Y_2^\top \\ Y_2 Y_1^\top & Y_2 Y_2^\top + \widetilde{W} \end{bmatrix} \right) \leq \lambda_{\min}(\widetilde{W}).$$

We are now ready to prove Theorem 18. Note that this result is very similar to that of [BHSW20], except that we work with the inverse trace rather than the minimum eigenvalue.

Proof. [Proof of Theorem 18] Let $\delta > 0$ be chosen later. We first argue that $\widehat{\text{tr}}$ must not be too large. Applying Proposition 24, we conclude that there is a universal constant $C' > 0$ such that $\text{tr}(W^{-1}) \leq C'd^2$ with probability at least $1/2$. Hence,

$$\begin{aligned} \mathbb{P}\{\widehat{\text{tr}} \leq CC'd^2\} &\geq \mathbb{P}\{\text{tr}(W^{-1}) \leq C'd^2 \text{ and } \widehat{\text{tr}} \leq C \text{tr}(W^{-1})\} \\ &\geq \mathbb{P}\{\text{tr}(W^{-1}) \leq C'd^2\} - \mathbb{P}\{\widehat{\text{tr}} > C \text{tr}(W^{-1})\} \geq \frac{1}{2} - \delta. \end{aligned}$$

Next, suppose for the sake of contradiction that $n \leq d/2$. Let \mathcal{F}_n denote the σ -algebra generated by the information available to the algorithm up to iteration n , that is, the queries v_1, \dots, v_n , the responses w_1, \dots, w_n , and any external randomness used by the algorithm (which is independent of W). Applying Propositions 21 and 22,

$$\begin{aligned} \mathbb{P}\{\widehat{\text{tr}} < C^{-1} \text{tr}(W^{-1})\} &\geq \mathbb{P}\{\widehat{\text{tr}} \leq CC'd^2 \text{ and } \lambda_{\max}(W^{-1}) > C^2 C' d^2\} \\ &\geq \mathbb{P}\{\widehat{\text{tr}} \leq CC'd^2 \text{ and } \lambda_{\max}(\widetilde{W}^{-1}) > C^2 C' d^2\} \\ &= \mathbb{E}[\mathbb{1}\{\widehat{\text{tr}} \leq CC'd^2\} \mathbb{P}\{\lambda_{\max}(\widetilde{W}^{-1}) \geq C^2 C' d^2 \mid \mathcal{F}_n\}]. \end{aligned}$$

According to Proposition 21, conditionally on \mathcal{F}_n , \widetilde{W} has the $\text{Wishart}(d-n)$ distribution. By applying Proposition 23,

$$\begin{aligned} \mathbb{P}\{\lambda_{\max}(\widetilde{W}^{-1}) \geq C^2 C' d^2 \mid \mathcal{F}_n\} &\geq \mathbb{P}\{\lambda_{\max}(\widetilde{W}^{-1}) \geq 4C^2 C' (d-n)^2 \mid \mathcal{F}_n\} \\ &= \mathbb{P}\left\{\lambda_{\min}(\widetilde{W}) \leq \frac{1}{4C^2 C' (d-n)^2} \mid \mathcal{F}_n\right\} \gtrsim \frac{1}{C\sqrt{C'}}. \end{aligned}$$

Therefore,

$$\mathbb{P}\{\widehat{\text{tr}} < C^{-1} \text{tr}(W^{-1})\} \gtrsim \mathbb{P}\{\widehat{\text{tr}} \leq CC'd^2\} \frac{1}{C\sqrt{C'}} \geq \frac{1/2 - \delta}{C\sqrt{C'}},$$

which is larger than δ provided that δ is chosen sufficiently small (depending only on C). This contradicts the success probability of the algorithm, and hence we deduce that $n \geq d/2$. \square

4.3 Useful facts about Wishart matrices

We collect together useful facts about Wishart matrices which are used in the proofs.

Proposition 23 (extreme singular values of a Gaussian matrix). *Let $W \sim \text{Wishart}(d)$. For any $x \in [0, 1]$,*

$$\mathbb{P}\{\lambda_{\min}(W) \leq \frac{x}{d^2}\} \asymp \sqrt{x}.$$

Also, there is a universal constant $C > 0$ such that

$$\mathbb{P}\{\lambda_{\max}(W) \geq C(1+t)\} \leq 2\exp(-dt).$$

Proof. See, e.g., [Ede89, Theorem 5.1] and [Ver18, Theorem 4.4.5]. \square

Proposition 24 (bound on the inverse trace). *Let $W \sim \text{Wishart}(d)$. Then, for any $\delta > 0$, with probability at least $1 - \delta$, it holds that $\text{tr}(W^{-1}) \leq C_\delta d^2$ where C_δ is a constant depending only on δ .*

Proof. According to [Sza91, Theorem 1.2], there is a universal constant $C > 0$ such that for each $j = 1, \dots, d$ and $\alpha \geq 0$,

$$\mathbb{P}\left\{\frac{1}{\lambda_j(W)} \geq \frac{d^2}{\alpha^2 j^2}\right\} \leq (C\alpha)^{j^2}.$$

Let $\alpha < 1/C$ and let $E_\alpha := \{1/\lambda_j(W) \geq d^2/(\alpha^2 j^2) \text{ for some } j = 1, \dots, d\}$. By the union bound,

$$\mathbb{P}(E_\alpha) \leq \sum_{j=1}^d (C\alpha)^{j^2} \lesssim \frac{1}{\sqrt{\log(1/(C\alpha))}}.$$

On the event E_α^c ,

$$\text{tr}(W^{-1}) \leq \sum_{j=1}^d \frac{d^2}{\alpha^2 j^2} = \frac{\pi^2 d^2}{6\alpha^2},$$

which is the claimed result upon taking α sufficiently small. \square

Remark. The proof only shows that $\mathbb{P}\{\text{tr}(W^{-1}) \geq \eta d^2\} \lesssim 1/\sqrt{\log \eta}$ for $\eta \gg 1$, which is not enough to conclude that $\mathbb{E} \text{tr}(W^{-1})$ is finite. In fact, it holds that $\mathbb{E} \text{tr}(W^{-1}) = \infty$, which can already be seen from Proposition 23.

5 A lower bound for sampling from Gaussians via reduction to block Krylov

In this section, we prove Theorem 3. Our proof proceeds in two parts: we first show a lower bound against the block Krylov method, and then a reduction showing that an arbitrary adaptive algorithm can be simulated via a block Krylov method.

5.1 Preliminaries

We first record some important facts that we will use later on. Throughout, let K be an odd integer. The following is a standard approximation-theoretic result:

Proposition 25 ([SV14, Proposition 2.4, rephrased]). *Let T_K be the degree- K Chebyshev polynomial, and let $1 = \beta_1 > \dots > \beta_{K+1} = -1$ be the set of real values β such that $T_K(\beta) \in \{-1, 1\}$. Then, for any real degree- K polynomial p such that $|p(\beta_i)| \leq 1$ for all β_i , we have $|p(x)| \leq |T_K(x)| \leq (|x| + \sqrt{x^2 - 1})^K$ for all $|x| > 1$.*

Let $c_0 > 0$ be a constant to be chosen later. The above proposition immediately implies:

Corollary 26 (approximation error). *Suppose that $K \leq c_0 \sqrt{\kappa} \log d$. Then, there exist $\kappa = \lambda_1 > \dots > \lambda_{K+2} = 1$ (that only depend on K and κ) such that for any real degree- K polynomial P , $\max_{1 \leq i \leq K+2} |\frac{1}{\lambda_i} - P(\lambda_i)| \geq d^{-2c_0 - O(1/\sqrt{\kappa})} / \kappa$.*

Proof. Set $\beta_1, \dots, \beta_{K+2}$ to be the solutions of $T_{K+1} \in \{-1, 1\}$, and for each $1 \leq i \leq K+2$, set $\lambda_i := \frac{(\kappa-1)}{2}(\beta_i + 1) + 1$; by construction, $\kappa = \lambda_1 > \dots > \lambda_{K+2} = 1$. Given any polynomial Q of degree at most $K+1$, note that if $|Q(\lambda_i)| \leq 1$ for all i , then the polynomial p given by $p(x) := Q(\frac{\kappa-1}{2}(x+1) + 1)$ satisfies $|p(\beta_i)| \leq 1$ for all i . By Proposition 25, for $x_0 := -(1 + \frac{2}{\kappa-1})$,

$$\begin{aligned} |Q(0)| = |p(x_0)| &\leq (|x_0| + \sqrt{x_0^2 - 1})^{K+1} \leq \left(1 + \frac{2}{\sqrt{\kappa}} + O\left(\frac{1}{\kappa}\right)\right)^{K+1} \\ &< \exp\left(\left(\frac{2}{\sqrt{\kappa}} + O\left(\frac{1}{\kappa}\right)\right)(c_0 \sqrt{\kappa} \log d + 1)\right) = d^{2c_0 + O(1/\sqrt{\kappa})}. \end{aligned}$$

Next, for a degree- K polynomial P , consider $Q(x) := d^{2c_0 + O(1/\sqrt{\kappa})}(1 - xP(x))$. Note that Q has degree $K+1$ and $|Q(0)| = d^{2c_0 + O(1/\sqrt{\kappa})}$, which implies that $|Q(\lambda_i)| > 1$ for some i , which in turn shows that $|\frac{1}{\lambda_i} - P(\lambda_i)| \geq d^{-2c_0 - O(1/\sqrt{\kappa})} / \kappa$. \square

We also introduce random matrix ensembles that are used in the proof, together with basic facts and properties.

Interestingly, as in the previous section, Wishart matrices are also useful for understanding block Krylov algorithms, but for a completely different reason. This time, we will study inner products between random vectors, which is also captured by a Wishart matrix. We denote by $\text{Wishart}(K, N)$ the law of the random matrix $XX^\top \in \mathbb{R}^{K \times K}$, where the entries of $X \in \mathbb{R}^{K \times N}$ are i.i.d. *standard* Gaussians. Note that this is a different convention from the previous section, in which each entry of X was i.i.d. $\mathcal{N}(0, \frac{1}{d})$.

We also define the Gaussian orthogonal ensemble (GOE) of size K , denoted $\text{GOE}(K)$. This is the law of a random symmetric matrix $G \in \mathbb{R}^{K \times K}$ where each diagonal entry $G_{i,i}$ is distributed as $\mathcal{N}(0, 1)$, and each off-diagonal entry $G_{i,j} = G_{j,i}$ is distributed as $\mathcal{N}(0, \frac{1}{2})$. Also, the entries $\{G_{i,j} : 1 \leq i \leq K, j \leq i\}$ are independent.

A long line of work (see, e.g., [JL15; BDER16; BG18; RR19; BBH21; Mik22]) shows that when $N \gg K^3$, the Wishart ensemble is well-approximated by a scaled and shifted GOE, a fact which we shall invoke in the sequel.

Lemma 27 (equivalence of Wishart and GOE). *Let $W \sim \text{Wishart}(K, N)$ be drawn from the Wishart distribution, and let W_0 be drawn from the distribution of symmetric matrices where the diagonal and above-diagonal entries are mutually independent, each diagonal entry is drawn as $\mathcal{N}(N, 2N)$,*

and each above-diagonal entry is drawn as $\mathcal{N}(0, N)$. (Equivalently, we can write $W_0 = NI + \sqrt{2N} G$, where $G \sim \text{GOE}(K)$.) Then,

$$\|\text{law}(W) - \text{law}(W_0)\|_{\text{TV}} \leq O\left(\frac{K^{3/2}}{N^{1/2}}\right).$$

Finally, we also require the following basic linear algebraic fact:

Proposition 28 (rotating the right singular vectors). *Let $V, V' \in \mathbb{R}^{K \times N}$ be such that $VV^\top = (V')(V')^\top$. Then, there exists an orthogonal matrix $U \in \mathbb{R}^{N \times N}$ such that $VU = V'$.*

5.2 Lower bound against block Krylov algorithms

We start with the following proposition, which will be useful in establishing the existence of matrices with different inverse traces but which generate similar power method iterates.

Proposition 29 (polynomial approximation and duality). *Suppose that $K \leq c_0 \sqrt{\kappa} \log d$. Then, there exist $\kappa = \lambda_1 > \lambda_2 > \dots > \lambda_{K+2} = 1$ and non-negative real numbers $x_1, \dots, x_{K+2}; x'_1, \dots, x'_{K+2}$, such that:*

1. For all $0 \leq j \leq K$, $\sum_{i=1}^{K+2} x_i \lambda_i^j = \sum_{i=1}^{K+2} x'_i \lambda_i^j$.
2. $\sum_{i=1}^{K+2} x_i = \sum_{i=1}^{K+2} x'_i = d$.
3. $\sum_{i=1}^{K+2} x_i / \lambda_i - \sum_{i=1}^{K+2} x'_i / \lambda_i \geq 2d^{1-2c_0-O(1/\sqrt{\kappa})} / \kappa$.

Proof. If we fix the values of the λ_i to be the choices in Corollary 26, this becomes a linear program in the variables $\{x_i\}_{i=1}^{K+2}, \{x'_i\}_{i=1}^{K+2}$. By writing $\mathbf{x} = (x_1, \dots, x_{K+2}, x'_1, \dots, x'_{K+2})$, our goal is to maximize $\mathbf{c}^\top \mathbf{x}$ over $\mathbf{x} \geq 0$ subject to $A\mathbf{x} = \mathbf{b}$. In our case, we set

$$\mathbf{c} := \begin{bmatrix} \lambda_1^{-1} \\ \vdots \\ \lambda_{K+2}^{-1} \\ -\lambda_1^{-1} \\ \vdots \\ -\lambda_{K+2}^{-1} \end{bmatrix}, \quad A := \begin{bmatrix} 1 & \dots & 1 & 1 & \dots & 1 \\ 1 & \dots & 1 & -1 & \dots & -1 \\ \lambda_1 & \dots & \lambda_{K+2} & -\lambda_1 & \dots & -\lambda_{K+2} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \lambda_1^K & \dots & \lambda_{K+2}^K & -\lambda_1^K & \dots & -\lambda_{K+2}^K \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 2d \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

We can consider the dual linear program, and by strong duality this maximization is equivalent to minimizing $\mathbf{b}^\top \mathbf{y}$ over \mathbf{y} such that $A^\top \mathbf{y} \geq \mathbf{c}$. By writing $\mathbf{y} = (z, y_0, y_1, \dots, y_K)$, this means we wish to minimize $2dz$ subject to $z + (y_0 + y_1 \lambda_i + \dots + y_K \lambda_i^K) \geq \frac{1}{\lambda_i}$ and $z - (y_0 + y_1 \lambda_i + \dots + y_K \lambda_i^K) \geq -\frac{1}{\lambda_i}$ for all $1 \leq i \leq K+2$. Equivalently, we wish to minimize $2dz$ subject to the existence of a polynomial P of degree at most K (with coefficients y_0, \dots, y_K) such that $z \geq |\frac{1}{\lambda_i} - P(\lambda_i)|$ for all $i \leq K+2$.

The minimum for the dual linear program (and thus the maximum for the primal linear program), is $2d \inf_{P \in \mathcal{P}_K} \max_{1 \leq i \leq K+2} |\frac{1}{\lambda_i} - P(\lambda_i)|$, where \mathcal{P}_K is the set of polynomials of degree at most K with real coefficients. By Corollary 26, this quantity is at least $2d^{1-2c_0-O(1/\sqrt{\kappa})} / \kappa$. \square

We note that a slightly strengthened version of Proposition 29 holds. Let $0 < c_1 < 1$.

Corollary 30 (existence of good solutions). *Proposition 29 holds, where we also ensure that each x_i and x'_i is at least $\frac{d}{2(K+2)}$ and $\frac{|x_i - x'_i|}{x_i} \leq \frac{2c_1}{1-c_1}$, though the right-hand side of the third condition becomes $c_1 d^{1-2c_0-O(1/\sqrt{\kappa})} / \kappa$.*

Proof. First, replace every x_i with $\frac{1}{2}(x_i + \frac{d}{K+2})$ and x'_i with $\frac{1}{2}(x'_i + \frac{d}{K+2})$. Then, we have that the replaced x_i, x'_i are at least $\frac{d}{2(K+2)}$, and the remaining statements in Proposition 29 hold, except the third which has the right-hand side replaced with $d^{1-2c_0-O(1/\sqrt{\kappa})}/\kappa$.

Next, replace every x_i with $\tilde{x}_i := \frac{1+c_1}{2}x_i + \frac{1-c_1}{2}x'_i$, and every x'_i with $\tilde{x}'_i := \frac{1+c_1}{2}x'_i + \frac{1-c_1}{2}x_i$. We still have that every $\tilde{x}_i, \tilde{x}'_i$ is at least $\frac{d}{2(K+2)}$, the first two conditions still hold, and the right-hand side of third condition is now $c_1 d^{1-2c_0-O(1/\sqrt{\kappa})}/\kappa$. Finally, note that $|\tilde{x}_i - \tilde{x}'_i| \leq c_1 |x_i - x'_i|$, whereas $\tilde{x}_i \geq \frac{1-c_1}{2}(x_i + x'_i)$. This implies that $\frac{|\tilde{x}_i - \tilde{x}'_i|}{\tilde{x}_i} \leq \frac{2c_1}{1-c_1}$. \square

We now have the necessary tools to prove our lower bound against block Krylov algorithms. Before doing so, we establish that there exist diagonal matrices D, D' which have substantially different inverse traces, but block Krylov algorithms cannot distinguish between them. To prove our actual lower bound, we show the same claim holds even if D, D' are randomly rotated, and the inverse trace difference is enough for a single sample to distinguish between them.

Lemma 31 (construction of diagonal matrices). *Suppose that $K \leq c_0 \sqrt{\kappa} \log d$ and $K \leq O(d)$. Then, there exist diagonal matrices $D, D' \in \mathbb{R}^{d \times d}$ with all diagonal entries between 1 and κ with the following properties.*

1. $|\text{tr}(D^{-1}) - \text{tr}(D'^{-1})| \geq c_1 d^{1-2c_0-O(1/\sqrt{\kappa})}/\kappa - 2(K+2)$.
2. Consider sampling K d -dimensional random vectors $v^{(1)}, \dots, v^{(K)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I_d)$. Then, the distributions of $\{\langle v^{(k)}, D^j v^{(\ell)} \rangle\}_{j \leq K+2; k, \ell \leq K}$ and $\{\langle v^{(k)}, D'^j v^{(\ell)} \rangle\}_{j \leq K+2; k, \ell \leq K}$ differ in total variation distance by at most $O(c_1 K^3 + K^3/d^{1/2})$.

Proof. Choose $\{x_i\}_{i=1}^{K+2}, \{x'_i\}_{i=1}^{K+2}$, and $\{\lambda_i\}_{i=1}^{K+2}$ satisfying Corollary 30. Define integers $\{N_i\}_{i=1}^{K+2}$ such that each N_i is either $\lfloor x_i \rfloor$ or $\lceil x_i \rceil$ and $\sum_{i=1}^{K+2} N_i = d$; define $\{N'_i\}_{i=1}^{K+2}$ similarly in terms of $\{x'_i\}_{i=1}^{K+2}$. We let D, D' be diagonal matrices such that for all i , D has N_i diagonal entries equal λ_i , and D' has N'_i diagonal entries equal to λ_i . Now, let $v^{(1)}, \dots, v^{(K)} \in \mathbb{R}^d$ be K random vectors drawn i.i.d. from $\mathcal{N}(0, I_d)$ (and define $v^{(1)'}, \dots, v^{(K)'} \in \mathbb{R}^d$ similarly). For $1 \leq i \leq K+2$, we define $v^{(k,i)}$ to be the projection of $v^{(k)}$ onto the dimensions corresponding to the diagonal entry λ_i for D . Note that $\{v^{(k,i)}\}_{i \leq K+2, k \leq K}$ are independent, and $v^{(k,i)} \sim \mathcal{N}(0, I_{N_i})$. Likewise, define $\{v^{(k,i)'}\}_{i \leq K+2, k \leq K}$ accordingly in terms of D' .

Note that $\text{tr}(D^{-1}) - \text{tr}(D'^{-1}) = \sum_{i=1}^{K+2} N_i/\lambda_i - \sum_{i=1}^{K+2} N'_i/\lambda_i$. Since $|N_i - x_i|, |N'_i - x'_i| \leq 1$, and since each $\lambda_i \geq 1$, it implies

$$\text{tr}(D^{-1}) - \text{tr}(D'^{-1}) \geq \frac{c_1 d^{1-2c_0-O(1/\sqrt{\kappa})}}{\kappa} - 2(K+2).$$

Next, we let $W^{(i)}$ represent the $K \times K$ matrix with entries $W_{k,\ell}^{(i)} = \langle v^{(k,i)}, v^{(\ell,i)} \rangle$ and define $W^{(i)'} \in \mathbb{R}^{K \times K}$ similarly. Note that the matrices $W^{(i)}, W^{(i)'}$ over all i are independent. In addition, $W^{(i)}$ has the $\text{Wishart}(K, N_i)$ distribution, and $W^{(i)'}$ has the $\text{Wishart}(K, N'_i)$ distribution. In addition, for any $k, \ell \leq K$ and $j \leq T$, we have that $\langle v^{(k)}, D^j v^{(\ell)} \rangle = \sum_{i=1}^{K+2} \lambda_i^j W_{k,\ell}^{(i)}$.

Now, we attempt to design a coupling between the matrices $\{W^{(i)}\}_{i=1}^{K+2}$ and $\{W^{(i)'}\}_{i=1}^{K+2}$ such that $W^{(i)} - W^{(i)'} = (x_i - x'_i) I_K$ for all $i \leq K+2$, with high probability. Note that this implies our claim, due to Corollary 30. To design this coupling, first note that by Lemma 27, if we draw $Z^{(i)} \sim N_i I_K + \sqrt{2N_i} \text{GOE}(K)$, then $\|\text{law}(W^{(i)}) - \text{law}(Z^{(i)})\|_{\text{TV}} \leq O(K^{3/2}/N_i^{1/2})$, and a similar statement holds if we define $Z^{(i)'}$ and compare its law to that of $W^{(i)'}$.

Note that the entries of $Z^{(i)}$ and $Z^{(i)'}$ are independent (apart from the requirement of symmetry), so we will attempt a coupling between the entries $Z_{k,\ell}^{(i)}$ and $Z_{k,\ell}^{(i)'}$. For $k < \ell$, since $Z_{k,\ell}^{(i)} \sim \mathcal{N}(0, N_i)$

and $Z_{k,\ell}^{(i)'} \sim \mathcal{N}(0, N_i')$, the total variation distance between their distributions is bounded up to a constant, using Corollary 30, by

$$\left| \frac{N_i'}{N_i} - 1 \right| \leq \left| \frac{x_i'}{x_i} - 1 \right| + \left| \frac{N_i' - x_i'}{N_i} \right| + \left| \frac{x_i' (N_i - x_i)}{N_i x_i} \right| \leq O\left(c_1 + \frac{K}{d}\right)$$

under our assumptions. Therefore, we can couple $Z_{k,\ell}^{(i)}$ and $Z_{k,\ell}^{(i)'}$ such that they fail to coincide with this probability. For $k = \ell$, we have $Z_{k,k}^{(i)} \sim \mathcal{N}(N_i, 2N_i)$ and $Z_{k,k}^{(i)'} + x_i - x_i' \sim \mathcal{N}(N_i' + x_i - x_i', 2N_i')$. The total variation distance between their distributions is bounded by a constant times

$$\left| \frac{N_i'}{N_i} - 1 \right| + \frac{|N_i' - x_i' + x_i - N_i|}{\sqrt{N_i}} \leq O\left(c_1 + \frac{K^{1/2}}{d^{1/2}}\right).$$

Therefore, we can couple the two random variables together so that $Z_{k,k}^{(i)} = Z_{k,k}^{(i)'} + x_i - x_i'$ fails with the above probability.

By a union bound, the coupling $Z^{(i)} = Z^{(i)'} + (x_i - x_i') I_K$ for all i fails with probability at most

$$O\left(K^3 \left(c_1 + \frac{K}{d}\right) + K^2 \left(c_1 + \frac{K^{1/2}}{d^{1/2}}\right)\right) = O\left(c_1 K^3 + \frac{K^{5/2}}{d^{1/2}}\right).$$

We dropped the $c_1 K^4/d$ term because of our assumption $K \leq O(d)$. Combining this with comparison between the Wishart and GOE ensembles and another union bound, we obtain the result. \square

Finally, we are able to prove our main lower bound against block Krylov algorithms.

Lemma 32 (lower bound against block Krylov algorithms). *Let κ, K, D, D' be as in Lemma 31. Then, let U be a uniformly random orthogonal matrix in $\mathbb{R}^{d \times d}$, and let $\Lambda = U^\top D U$ and $\Lambda' = U^\top D' U$. Let $v^{(1)}, \dots, v^{(K)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I_d)$. Then, for any $\delta > 0$, provided $K \leq O_\delta(\sqrt{\kappa} \log d)$ and $\kappa \leq d^{1/5-\delta}$, the distributions of $\{\Lambda^{j,v^{(k)}}\}_{j \leq (K+2)/2, k \leq K}$ and $\{\Lambda'^{j,v^{(k)}}\}_{j \leq (K+2)/2, k \leq K}$ differ in total variation distance by at most $o(1)$. On the other hand, drawing a sample either from $\mathcal{N}(0, \Lambda^{-1})$ or $\mathcal{N}(0, \Lambda'^{-1})$ can, with probability $1 - o(1)$, distinguish between the two cases.*

Proof. The following calculations are contingent on the values of the various parameters that we will choose at the end of the proof. From Lemma 31, there is a coupling such that the tuples $\{\langle v^{(k)}, D^j v^{(\ell)} \rangle\}_{j \leq K+2, k, \ell \leq K}$ and $\{\langle v^{(k)'}, D'^j v^{(\ell)'} \rangle\}_{j \leq K+2, k, \ell \leq K}$ are equal with high probability. In particular, it holds that $\langle D^i v^{(k)}, D^j v^{(\ell)} \rangle = \langle D'^i v^{(k)'}, D'^j v^{(\ell)'} \rangle$ for all $i, j \leq (K+2)/2$ and $k \leq K$ with high probability. By Proposition 28, there is a unitary matrix U_0 such that $D'^j v^{(k)'} = U_0 D^j v^{(k)}$ for all $j \leq (K+2)/2$ and all $k \leq K$ with high probability. Note then that the tuples $\{U^\top D^j U U^\top v^{(k)}\}_{j \leq (K+2)/2, k \leq K}$ and $\{U^\top U_0^\top D'^j U_0 U U^\top U_0^\top v^{(k)'}\}_{j \leq (K+2)/2, k \leq K}$ are equal with high probability, and this is a coupling which witnesses the fact that the distributions of $\{\Lambda^{j,v^{(k)}}\}_{j \leq (K+2)/2, k \leq K}$ and $\{\Lambda'^{j,v^{(k)'}}\}_{j \leq (K+2)/2, k \leq K}$ are at most $O(c_1 K^3 + K^3/d^{1/2})$ apart in total variation distance.

Finally, we note that from a single sample it is easy to distinguish between $\mathcal{N}(0, \Lambda^{-1})$ and $\mathcal{N}(0, \Lambda'^{-1})$. This is because if $X \sim \mathcal{N}(0, \Lambda^{-1})$, then $\mathbb{E}[\|X\|^2] = \text{tr}(\Lambda^{-1}) = \text{tr}(D^{-1}) = \sum_{i=1}^{K+2} N_i/\lambda_i$, but one checks that $\text{var}(\|X\|^2) = O(\sum_{i=1}^{K+2} N_i/\lambda_i^2) \leq O(d)$. Likewise, if $X' \sim \mathcal{N}(0, \Lambda'^{-1})$, then we have $\mathbb{E}[\|X'\|^2] = \sum_{i=1}^{K+2} N_i'/\lambda_i$ but $\text{var}(\|X'\|^2) = O(d)$. So, the difference in their expectations at least $c_1 d^{1-2c_0-O(1/\sqrt{\kappa})}/\kappa - 2(K+2)$, whereas the standard deviations are bounded by $O(d^{1/2})$.

To finish the proof, we must choose the values of c_0 and c_1 . We require the following conditions:

1. $c_1 K^3 = o(1)$.

2. $K^3/d^{1/2} = o(1)$.
3. $d^{1/2} = o(c_1 d^{1-2c_0-O(1/\sqrt{\kappa})}/\kappa - 2(K+2))$.

For the second condition, we can assume $\kappa \leq d^{1/3}/\log^4(d)$. To satisfy the first condition, we can set $c_1 = 1/(\kappa^{3/2} \log^4(d))$. Finally, if κ is sufficiently large and if c_0 is chosen depending on δ , then the third condition requires $\sqrt{\kappa} \log d + d^{1/2} = o(d^{1-\delta}/\kappa^{5/2})$, and it suffices for $\kappa \leq d^{1/5-\delta}$. \square

Remark. We did not attempt to optimize the exponent in the condition $\kappa \leq d^{1/5-\delta}$. Indeed, by using the chain rule for the KL divergence rather than a union bound in the proof of Lemma 31, we believe that the total variation bound can be improved to $O(c_1 K^{3/2} + K^{5/2}/d^{1/2})$, and a back-of-the-envelope calculation suggests that this could improve the condition to $\kappa \leq d^{2/7-\delta}$. Nevertheless, this falls short of capturing the full regime $\sqrt{\kappa} \log d \leq O(d)$, and we leave this as an open question.

5.3 Reduction to block Krylov algorithms

In this section, we show that in order to prove a lower bound for sampling from Gaussians against any query algorithm, it suffices to prove a lower bound against block Krylov algorithms.

5.3.1 Setup

Let $\Lambda = U^\top D U$, where D is a (possibly random) diagonal matrix, U is a Haar-random orthogonal matrix, and U and D are independent. We consider the following model, which is a strengthening of the matrix-vector product model:

Definition 33 (extended oracle model). *Given $K \in \mathbb{N}$, for all $k \in [K]$, the algorithm chooses a new query point v_k , and receives the information $\{\Lambda^i v_j\}_{(i,j) \in H_k}$, where $H_k := \{(i,j) : i+j \leq k+1, i \geq 0, 1 \leq j \leq k\}$ is a set of ordered pairs of nonnegative integers. We use the following notation $\{\Lambda^i v_j\}_S$ for any set S to denote $\{\Lambda^i v_j\}_{(i,j) \in S}$.*

This is clearly a stronger oracle model than before, so a lower bound against algorithms in the extended oracle model implies a lower bound against algorithms in the original matrix-vector model.

Definition 34 (adaptive deterministic algorithm). *An adaptive deterministic algorithm \mathcal{A} that makes K extended oracle queries (see Definition 33) is given by a deterministic collection of functions $v_1, v_2(\cdot), \dots, v_K(\cdot)$, where v_1 is constant and each $v_k(\cdot)$ is a function of $\frac{k(k+1)}{2} - 1$ inputs. This corresponds to a sequence of queries where the k -th query $v_k(\{\Lambda^i v_j\}_{H_{k-1}})$ is chosen adaptively based on the information available to the algorithm at the start of iteration k . (Note that v_1 has no inputs.) When the choice of the inputs is clear from context, we may simply write $v_k = v_k(\{\Lambda^i v_j\}_{H_{k-1}})$.*

In the extended oracle model, the next lemma shows that we can assume that each v_k is a unit vector orthogonal to its inputs.

Lemma 35 (extended oracle and orthogonal queries). *For $k \in [2, K]$, let v_k be as stated in Definition 34 and let $\{\Lambda^i v_j\}_{H_{k-1}}$ be as stated in Definition 33. Then, without loss of generality, we may assume that v_k is orthogonal to the subspace spanned by the vectors in $\{\Lambda^i v_j\}_{H_{k-1}}$.*

Proof. Assume for sake of contradiction that this were not the case. Then, we can decompose $v_k = \sum_{(i,j) \in H_{k-1}} c_{i,j} \Lambda^i v_j + c^\perp v_k^\perp$ where v_k^\perp is a unit vector orthogonal to $\{\Lambda^i v_j\}_{H_{k-1}}$ and each $c_{i,j}$ and c^\perp is a scalar. At the end of iteration k , the new information obtained by the algorithm is $\{\Lambda^i v_j\}_{i+j=k+1, j \leq k}$. For all $(i,j) \neq (1,k)$, the new information does not depend on v_k . Also,

$\Lambda v_k = \sum_{(i,j) \in H_{k-1}} c_{i,j} \Lambda^{i+1} v_j + c^\perp \Lambda v_k^\perp$, where each $\Lambda^{i+1} v_j$ is information obtained by the algorithm at the end of iteration $k+1$ regardless (due to our extended query model). Since $(i+1, j) \in H_k$ if $(i, j) \in H_{k-1}$, and since $(1, k) \in H_k$, this expression shows that the algorithm would receive the same amount of information (or more, if $c^\perp = 0$) if it queries v_k^\perp instead of v_k . Applying this reasoning inductively proves the claim. \square

We compare to a *block Krylov algorithm*, which makes i.i.d. standard Gaussian queries z_1, \dots, z_K and then receives $\{\Lambda^i z_j\}$ for all $i, j \leq K$. Recall that a block Krylov algorithm does not make *adaptive* queries, so it is easier to prove lower bounds against block Krylov algorithms. Our goal is to now show that block Krylov algorithms can simulate an adaptive deterministic algorithm.

5.3.2 Conditioning lemma

We start by proving a general conditioning lemma which will be invoked repeatedly in the reduction to block Krylov algorithms. This lemma roughly shows that if the adaptive algorithm knows $\{\Lambda^i v_j\}_{H_k}$, the posterior distribution of Λ given $\{\Lambda^i v_j\}_{H_k}$ is indeed rotationally symmetric on the orthogonal complement $\{\Lambda^i v_j\}_{H_k}^\perp$.

We will use the notation $\stackrel{d}{=}$ to denote that two random variables are equal in probability distribution (possibly conditioned on other information).

Lemma 36 (conditioning lemma, preliminary version). *Let U be a Haar-random orthogonal matrix, and $\Lambda = U^\top D U$, where D is a (possibly random) positive diagonal matrix. Suppose that \mathcal{A} is an adaptive deterministic algorithm that generates extended oracle queries v_1, \dots, v_K , and after the k -th query knows $\Lambda^i v_j$ for all $(i, j) \in H_k$. For any integer $m \geq 1$, let k be the integer such that $\frac{k(k+1)}{2} \leq m < \frac{(k+1)(k+2)}{2}$, i.e., m is at least the k -th triangular number but less than the $(k+1)$ -th triangular number. Consider the order of vectors $v_1, \Lambda v_1, v_2, \Lambda^2 v_1, \Lambda v_2, v_3, \Lambda^3 v_1, \dots$ (this enumerates $\Lambda^i v_j$ in order of $i+j$, breaking ties with smaller values of j first). Let W_m be the set of first m of these vectors and X_k be the set $\{v_1, \dots, v_k\}$. Let V be a Haar-random orthogonal matrix fixing W_m and acting on the orthogonal complement W_m^\perp . Then, $(X_k, U) \stackrel{d}{=} (X_k, UV)$.*

Before proving this lemma, we note that since the algorithm is deterministic and D is fixed, W_m and X_k are deterministic functions of Λ , and thus of U . Hence, we can write $v_k(U'), W_m(U'), X_k(U')$ to be the v_k, W_m, X_k that would have been generated if we started with $\Lambda' = (U')^\top \Lambda U'$. (If no argument is given, v_k, W_m, X_k are assumed to mean $v_k(U), W_m(U), X_k(U)$, respectively.) We note the following proposition.

Proposition 37 (fixing the first m queries and responses). *Suppose that V is any orthogonal matrix fixing $W_m(U)$. Then, $W_m(U) = W_m(UV)$.*

Proof. We prove $W_{m'}(U) = W_{m'}(UV)$ for all $m' \leq m$. The base case of $k = 1$ is trivial, since v_1 is fixed. We now prove the induction step for m' .

If $m' \leq m$ is a triangular number, $m' = \frac{k(k+1)}{2}$, then the m' -th vector in W_m is v_k . But note that $v_k(U)$ is a deterministic function of $W_{m'-1}(U)$, and $v_k(UV)$ is the same deterministic function of $W_{m'-1}(UV)$. Hence, if the induction hypothesis holds for $m' - 1$, it also holds for m .

If $m' \leq m$ is not a triangular number, then the m' -th number in $W_m(U)$ is $\Lambda^i v_j$ for some $i \geq 1$. Likewise, the m' -th number in $W_m(UV)$ is $V^\top \Lambda^i V v_j(UV)$. Since $i \geq 1$, we know that $v_j(U) = v_j(UV)$, by the induction hypothesis on $\frac{j(j+1)}{2} < m'$. But, we know that V fixes W_m , which means it fixes v_j and $\Lambda^i v_j$. Thus, $V^\top \Lambda^i V v_j(UV) = V^\top \Lambda^i V v_j = \Lambda^i v_j$. \square

We are now ready to prove Lemma 36.

Proof. [Proof of Lemma 36] We prove this by induction on m . For the base case $m = 1$, U is a random matrix and V is a random matrix that fixes v_1 . Note that v_1 is chosen independently of Λ (and thus of U), so U and V are independent. Even for any fixed V , the distribution UV is a uniformly random orthogonal matrix, so overall $U \stackrel{d}{=} UV$. Also, v_1 is deterministic, so $(v_1, U) \stackrel{d}{=} (v_1, UV)$.

For the induction step, we split the proof into 2 cases. The proofs in both cases will be very similar, but with minor differences.

Case 1: m is a triangular number. This means that the m -th vector added is v_k , where $m = \frac{k(k+1)}{2}$. Let V_1 be a random orthogonal matrix fixing W_{m-1} and V_2 be a random orthogonal matrix fixing W_m . Our goal is then to show $(X_k, U) \stackrel{d}{=} (X_k, UV_2)$.

To make this rigorous, we note an order of generating the random variables. First, we generate U randomly: W_m and X_k are deterministic in terms of U . Next, we define V_1 to be a random rotation fixing W_{m-1} . Finally, we define V_2 to be a random rotation fixing W_m , where V_1, V_2 are conditionally independent on U .

First, we prove that $(X_k, U) \stackrel{d}{=} (X_k, UV_1)$. Note that $U \stackrel{d}{=} UV_1$ by our inductive hypothesis. In addition, since V_1 fixes $W_{m-1}(U)$, $W_{m-1}(U) = W_{m-1}(UV_1)$ by Proposition 37. Since $m = \frac{k(k+1)}{2}$ is a triangular number, $X_k(\cdot)$ is a deterministic function of $W_{m-1}(\cdot)$, which means $X_k(U) = X_k(UV_1)$. Hence, $(X_k, U) \stackrel{d}{=} (X_k(UV_1), UV_1) = (X_k, UV_1)$.

Next, we prove that $(X_k, UV_2) \stackrel{d}{=} (X_k, UV_1V_2)$. It suffices to prove that

$$(X_k, U, V_2) \stackrel{d}{=} (X_k, UV_1, V_2).$$

To do so, we first show that $V_2 = f(U, R)$, where f is a deterministic function and R represents a random orthogonal matrix over $d - \dim(W_m)$ dimensions that is independent of U . (Recall that W_m is a deterministic function of U .) To define $f(U, R)$, we consider some deterministic map that sends each W_m to a set of $d - \dim(W_m)$ basis vectors in W_m^\perp . We then define $V_2 = f(U, R)$ to act on W_m^\perp using R and the correspondence of basis vectors. Since W_m and X_k are deterministic in terms of U , this means $f(U, R)$ is well-defined. We will now show that

$$V_2 = f(U, R) = f(UV_1, R) \quad \text{and} \quad X_k = X_k(UV_1).$$

Since $U \stackrel{d}{=} UV_1$ by our inductive hypothesis,

$$(X_k, U, V_2) \stackrel{d}{=} (X_k(UV_1), UV_1, f(UV_1, R)) = (X_k, UV_1, V_2).$$

By Proposition 37, $W_{m-1}(U) = W_{m-1}(UV_1)$, and since $X_k(\cdot)$ is deterministic given $W_{m-1}(\cdot)$ for $m = \frac{k(k+1)}{2}$, $X_k(U) = X_k(UV_1)$. This implies $W_m(U) = W_m(UV_1)$, which means $f(UV_1, R) = f(U, R)$, since $f(\cdot, R)$ only depends on $W_m(\cdot)$ and R . This completes the proof.

Next, we show that $(X_k, UV_1V_2) \stackrel{d}{=} (X_k, UV_1)$. Since we chose the order with U being defined first, we are allowed to condition on U . Since X_k is deterministic in terms of U , it suffices to show that $V_1V_2 \mid U \stackrel{d}{=} V_1 \mid U$. Since W_{m-1}, W_m are also deterministic given U , note that V_1 is a uniformly random orthogonal matrix fixing W_{m-1} , and V_2 is a random orthogonal matrix fixing $W_m \supset W_{m-1}$. Since V_1 and V_2 are conditionally independent given U , this means $V_1V_2 \mid U$ is a uniformly random orthogonal matrix fixing W_{m-1} , so $V_1V_2 \mid U \stackrel{d}{=} V_1 \mid U$.

In summary, we have that

$$\begin{aligned}(X_k, U) &\stackrel{d}{=} (X_k, UV_1) \\ &\stackrel{d}{=} (X_k, UV_1V_2) \\ &\stackrel{d}{=} (X_k, UV_2).\end{aligned}$$

Case 2: m is not a triangular number. Again, let V_1 be a random orthogonal matrix fixing W_{m-1} and V_2 be a random orthogonal matrix fixing W_m . Our goal is again to show that $(X_k, U) \stackrel{d}{=} (X_k, UV_2)$.

First, we again have $(X_k, UV_1) \stackrel{d}{=} (X_k, U)$ by our inductive hypothesis.

Next, we show that $(X_k, UV_2) \stackrel{d}{=} (X_k, UV_2V_1)$. It suffices to prove that

$$(X_k, U, V_2) \stackrel{d}{=} (X_k, UV_1, V_1^\top V_2 V_1),$$

since $(UV_1)(V_1^\top V_2 V_1) = UV_2V_1$. We recall the random variable R and use the same function $V_2 = f(U, R)$. Since we have already shown that $U \stackrel{d}{=} UV_1$, this implies that $(X_k, U, V_2) \stackrel{d}{=} (X_k(UV_1), UV_1, f(UV_1, R))$. Since m is not triangular, $X_k(\cdot)$ is contained in $W_{m-1}(\cdot)$, so by Proposition 37, $X_k(U) = X_k(UV_1)$. So, we have

$$(X_k, U, V_2) \stackrel{d}{=} (X_k(UV_1), UV_1, f(UV_1, R)) = (X_k, UV_1, f(UV_1, R)).$$

Now, if we fix U and V_1 , $W_{m-1}(UV_1) = W_{m-1}(U)$ by Proposition 37. However, since the m -th (i, j) pair has $i \geq 1$ when m is not triangular, the final vector in $W_m(UV_1)$ will be $V_1^\top \Lambda^i V_1 v_j = V_1^\top (\Lambda^i v_j)$. For fixed U, V_1 , $f(U, R)$ is a random rotation fixing W_{m-1} and $\Lambda^i v_j$, but $f(UV_1, R)$ is a random rotation fixing W_{m-1} and $V_1^\top (\Lambda^i v_j)$. Since V_1^\top fixes W_{m-1} by how we defined V_1 , this means that for fixed U, V_1 , $f(U, R)$ is a random rotation fixing W_m but $f(UV_1, R)$ is a random rotation fixing $V_1^\top W_m$. Therefore, conditioned on U, V_1 , $f(UV_1, R)$ has the same distribution as $V_1^\top f(U, R) V_1$. Since X_k is deterministic in terms of U , this means

$$(X_k, UV_1, f(UV_1, R)) \mid U, V_1 \stackrel{d}{=} (X_k, UV_1, V_1^\top f(U, R) V_1) \mid U, V_1.$$

We can remove the conditioning to establish that $(X_k, UV_1, f(UV_1, R)) \stackrel{d}{=} (X_k, UV_1, V_1^\top f(U, R) V_1) = (X_k, UV_1, V_1^\top V_2 V_1)$, which completes the proof.

Next, we show that $(X_k, UV_2V_1) \stackrel{d}{=} (X_k, UV_1)$. The proof is essentially the same as in the case when m is triangular. We again condition on U , and we have that $V_2V_1 \mid U \stackrel{d}{=} V_1 \mid U$ have the same distribution as uniform orthogonal matrices fixing $W_{m-1}(U)$. Since X_k is a deterministic function of U , this means $(X_k, UV_2V_1) \mid U \stackrel{d}{=} (X_k, UV_1) \mid U$, and removing the conditioning finishes the proof.

In summary,

$$\begin{aligned}(X_k, U) &\stackrel{d}{=} (X_k, UV_1) \\ &\stackrel{d}{=} (X_k, UV_2V_1) \\ &\stackrel{d}{=} (X_k, UV_2).\end{aligned}$$

□

We now prove our main conditioning lemma, which will be a modification of Lemma 36.

Lemma 38 (conditioning lemma). *Let all notation be as in Lemma 36, and let V_0 be a fixed orthogonal matrix fixing W_m . Importantly, V_0 is a deterministic function only depending on W_m (and not directly on U). Then, $(X_k, U) \stackrel{d}{=} (X_k, UV_0)$.*

Proof. First, note that since V_0 is a deterministic function of W_m , it is also a deterministic function of U . We can write $V_0(\cdot)$ as this function, and $V_0 = V_0(U)$.

Now, Lemma 36 proves that $(X_k, U) \stackrel{d}{=} (X_k, UV)$. Note that conditioned on U , V is a random matrix fixing W_m and V_0 is a fixed matrix fixing W_m , which means that $VV_0 \mid U \stackrel{d}{=} V \mid U$. Hence, $(X_k, UV) \stackrel{d}{=} (X_k, UVV_0)$. But from Proposition 37, $X_k(UV) = X_k(U)$ and $W_m(UV) = W_m(U)$, which means that $V_0(\cdot)$, which only depends on $W_m(\cdot)$, satisfies $V_0(UV) = V_0(U)$. Hence, because $U \stackrel{d}{=} UV$, we have $(X_k, UVV_0) = (X_k(UV), UV \cdot V_0(UV)) \stackrel{d}{=} (X_k(U), U \cdot V_0(U)) = (X_k, UV_0)$.

In summary, we have that $(X_k, U) \stackrel{d}{=} (X_k, UV) \stackrel{d}{=} (X_k, UVV_0) \stackrel{d}{=} (X_k, UV_0)$, which completes the proof. \square

5.3.3 From query algorithms to block Krylov algorithms

In this section, we carry out the high-level outline from Section 2.2.2. We aim to prove the following result, which implies that any adaptive deterministic algorithm in the extended oracle model can be simulated by rotating the output of a block Krylov algorithm.

Lemma 39 (reduction to block Krylov). *Suppose $\Lambda = U^\top DU$, where U is a Haar-random orthogonal matrix and D is a diagonal matrix drawn from some (possibly unknown) distribution. Let $v_1, v_2(\cdot), \dots, v_K(\cdot)$ be an adaptive deterministic algorithm that makes K orthonormal queries, where $K^2 < d$. Let $v_1^{\text{alg}}, v_2^{\text{alg}}, \dots, v_K^{\text{alg}}$ be recursively defined as follows: $v_1^{\text{alg}} = v_1$, and $v_k^{\text{alg}} = v_k(\{\Lambda^i v_j^{\text{alg}}\}_{H_{k-1}})$ for $k \geq 2$. Let z_1, \dots, z_K be i.i.d. standard Gaussian vectors. Then, from the collection $\{\Lambda^i z_j\}_{H_K}$ (without knowledge of D or Λ), we can construct a set of unit vectors $\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_K$, and a set of rotation matrices $U_1^{\text{sim}}, U_2^{\text{sim}}, \dots, U_K^{\text{sim}}$, where \tilde{v}_k and U_k^{sim} only depend on $\{\Lambda^i z_j\}_{H_{k-1}}$ and z_k , and such that*

$$\{(U_{1:K}^{\text{sim}})^\top \Lambda^i \tilde{v}_j\}_{H_K} \stackrel{d}{=} \{\Lambda^i v_j^{\text{alg}}\}_{H_K},$$

where $U_{1:K}^{\text{sim}} := U_1^{\text{sim}} \dots U_K^{\text{sim}}$, and the equivalence in distribution is over the randomness of Λ and $\{z_i\}_{i \leq K}$. Moreover, $\{\Lambda^i \tilde{v}_j\}_{H_K}$ is deterministically determined by $\{\Lambda^i z_j\}_{H_K}$.

Lemma 39 says that the knowledge of $\Lambda^i z_j$ alone is sufficient to reconstruct the distribution of any adaptive algorithm's queries and responses. The proof of the lemma requires introducing a hefty amount of notation, but we emphasize that it follows along the lines of Section 2.2.2.

First, we describe how to construct \tilde{v}_k . Let $\tilde{v}_1 = \frac{z_1}{\|z_1\|}$, and for $k \geq 2$, let \tilde{v}_k be the unit vector parallel to the component of z_k that is orthogonal to the span of $\{\Lambda^i z_j\}_{H_{k-1}}$. (With probability 1, this is well-defined.)

Because each \tilde{v}_k is a linear combination of $\{\Lambda^i z_j\}_{H_{k-1}}$ and z_k , we can construct the set $\{\Lambda^i \tilde{v}_j\}_{H_K}$ from the set $\{\Lambda^i z_j\}_{H_K}$.

We now construct the rotation matrices U_k^{sim} . First, we define matrix-valued functions $U_k(\cdot)$, for $k = 1, \dots, K$, as follows.

Definition 40 (rotations fixing previous queries and responses). *For $1 \leq k \leq K$, the function $U_k(\cdot)$ takes arguments $\{x_{i,j}\}_{H_{k-1}}$, y_k , z_k , where the vectors y_k and z_k have unit norm and are both orthogonal to the collection $\{x_{i,j}\}_{H_{k-1}}$.*

To define $U_1(\cdot)$: since H_0 is empty, the first function U_1 only takes arguments y_1, z_1 , and is such that $U_1(y_1, z_1)$ is a deterministic orthogonal matrix that satisfies $U_1(y_1, z_1)^\top y_1 = z_1$. Note that $U_1(\cdot)$ exists because y_1 and z_1 both have unit norm; for example, we can complete y_1 and z_1 to orthonormal bases (y_1, y_2, \dots, y_d) , (z_1, z_2, \dots, z_d) and take $U_1(y_1, z_1) = \sum_{i=1}^d y_i z_i^\top$.

To define $U_k(\cdot)$: $U_k(\{x_{i,j}\}_{H_{k-1}}, y_k, z_k)$ is a deterministic orthogonal matrix that satisfies

$$\begin{aligned} U_k^\top x_{i,j} &= x_{i,j}, & \text{for all } (i,j) \in H_{k-1}, \\ U_k^\top y_k &= z_k. \end{aligned} \tag{5.1}$$

Such a choice of U_k is always possible, because $k^2 < d$, and because y_k and z_k are orthogonal to $x_{i,j}$; for example, we can start with the identity matrix on the subspace spanned by $\{x_{i,j}\}_{H_{k-1}}$ and add to it a sum of outer products formed by completing y_k and z_k to two orthonormal bases of the orthogonal complement.

Next, we describe how to construct U_k^{sim} . We will define U_k^{sim} along with an auxiliary sequence $\{v_k^{\text{sim}}\}_{k=1,2,\dots,K-1}$.

Definition 41 (simulated sequences). We let $v_1^{\text{sim}} = v_1$, and $U_1^{\text{sim}} = U_1(\tilde{v}_1, v_1^{\text{sim}})$. For $k \geq 2$, v_k^{sim} and U_k^{sim} are defined recursively as follows:

$$\begin{aligned} v_k^{\text{sim}} &= v_k(\{(U_{1:(k-1)}^{\text{sim}})^\top \Lambda^i \tilde{v}_j\}_{H_{k-1}}) \\ U_k^{\text{sim}} &= U_k(\{(U_{1:(k-1)}^{\text{sim}})^\top \Lambda^i \tilde{v}_j\}_{H_{k-1}}, (U_{1:(k-1)}^{\text{sim}})^\top \tilde{v}_k, v_k^{\text{sim}}). \end{aligned}$$

Intuitively, one can think of v_k^{sim} as the k th vector the simulator thinks the algorithm is querying, and U_k^{sim} as a rotation that corresponds v_k^{sim} to the random unit vector known by block Krylov.

Proposition 42 (existence of rotations). Each U_k^{sim} is well-defined.

Proof. To show that this choice of U_k^{sim} is possible, we need to check that $(U_{1:(k-1)}^{\text{sim}})^\top \tilde{v}_k, v_k^{\text{sim}}$ both have unit norm and are orthogonal to the subspace S_k spanned by $(U_{1:(k-1)}^{\text{sim}})^\top \Lambda^i \tilde{v}_j$ for $(i,j) \in H_{k-1}$. They both have unit norm because \tilde{v}_k and v_k^{sim} are constructed to have unit norm, and inductively we can assume $U_{1:(k-1)}^{\text{sim}}$ is orthogonal. Note that v_k^{sim} is orthogonal to S_k by our assumption on the function $v_k(\cdot)$, and $(U_{1:(k-1)}^{\text{sim}})^\top \tilde{v}_k$ is also orthogonal to S_k because

$$\langle (U_{1:(k-1)}^{\text{sim}})^\top \Lambda^i \tilde{v}_j, (U_{1:(k-1)}^{\text{sim}})^\top \tilde{v}_k \rangle = \langle \Lambda^i \tilde{v}_j, \tilde{v}_k \rangle = 0,$$

where the second line follows from the definition of \tilde{v}_k . □

We summarize some additional properties of v_k^{sim} and U_k^{sim} in the following lemma.

Lemma 43 (properties of the simulated sequences). The variables U_k^{sim} and v_k^{sim} for $k = 1, \dots, K$ defined above satisfy the following properties:

(P1) v_k^{sim} depends only on $\{\Lambda^i \tilde{v}_j\}_{H_{k-1}}$, and U_k^{sim} depends only on $\{\Lambda^i \tilde{v}_j\}_{i+j \leq k}$.

(P2) For any $k \geq j$, we have

$$\tilde{v}_j = U_{1:k}^{\text{sim}} v_j^{\text{sim}}.$$

(P3) For $k \geq 2$, v_k^{sim} satisfies

$$v_k^{\text{sim}} = v_k(\{(U_{1:(k-1)}^{\text{sim}})^\top \Lambda^i U_{1:(k-1)}^{\text{sim}} v_j^{\text{sim}}\}_{H_{k-1}}).$$

(P4) For $k \geq 2$, U_k^{sim} satisfies

$$U_k^{\text{sim}} = U_k(\{(U_{1:(k-1)}^{\text{sim}})^\top \Lambda^i U_{1:(k-1)}^{\text{sim}} v_j^{\text{sim}}\}_{H_{k-1}}, (U_{1:(k-1)}^{\text{sim}})^\top \tilde{v}_k, v_k^{\text{sim}}).$$

Proof. (P1) is immediate from the definitions, since $\{(i, j) : i + j \leq k\} = H_{k-1} \cup \{(0, k)\}$.

To show (P2), note that the second property of the function U_k from (5.1) implies that

$$v_j^{\text{sim}} = (U_j^{\text{sim}})^\top (U_{1:(j-1)}^{\text{sim}})^\top \tilde{v}_j = (U_{1:j}^{\text{sim}})^\top \tilde{v}_j.$$

This proves (P2) for $k = j$. To prove (P2) for $k > j$, we use induction on k . If (P2) holds for $k - 1 \geq j$, then

$$(U_{1:k}^{\text{sim}})^\top \tilde{v}_j = (U_k^{\text{sim}})^\top (U_{1:(k-1)}^{\text{sim}})^\top \tilde{v}_j = (U_{1:(k-1)}^{\text{sim}})^\top \tilde{v}_j = v_j^{\text{sim}}.$$

Above, the middle equality holds by the first property of (5.1), since U_k^{sim} fixes $(U_{1:(k-1)}^{\text{sim}})^\top \tilde{v}_j$ because $j \leq k - 1$. The final equality holds by our inductive hypothesis. So, (P2) holds for k .

Finally, (P3) and (P4) then follow from (P2), since $k - 1 \geq j$ if $j \in H_{k-1}$. \square

We highlight the importance of (P2) for $k = K$, which roughly states that $(U_{1:K}^{\text{sim}})^\top$ actually sends each block Krylov-generated vector \tilde{v}_j to the simulated vector v_j^{sim} .

Before proving Lemma 39, we must make one more basic definition.

Definition 44 (queries and data). For $k \geq 2$, given the matrix Λ and a set $\{v_j\}_{1 \leq j \leq k-1}$, define \mathfrak{C}_k as the function that satisfies $\mathfrak{C}_k(\Lambda, \{v_j\}_{1 \leq j \leq k-1}) = \{\Lambda^i v_j\}_{H_{k-1}}$. In addition, define $\mathfrak{D}_k = v_k \circ \mathfrak{C}_k$.

We are now ready to prove Lemma 39. Although the proof is notationally burdensome, the message is that we can show the equality of distributions inductively by repeatedly invoking the conditioning lemma (Lemma 38), which is designed precisely for the present situation.

Proof. [Proof of Lemma 39] For $1 \leq k \leq K$, let $\Lambda_k := (U_{1:k}^{\text{sim}})^\top \Lambda U_{1:k}^{\text{sim}}$. Since we can write $(U_{1:k}^{\text{sim}})^\top \Lambda^i \tilde{v}_j = (U_{1:k}^{\text{sim}})^\top \Lambda^i (U_{1:k}^{\text{sim}}) v_j^{\text{sim}} = \Lambda_k^i v_j^{\text{sim}}$ for any $k \geq j$ by (P2) of Lemma 43, it suffices to inductively prove that for all $1 \leq k \leq K$,

$$(\Lambda_k, \{v_j^{\text{sim}}\}_{1 \leq j \leq k}) \stackrel{d}{=} (\Lambda, \{v_j^{\text{alg}}\}_{1 \leq j \leq k}).$$

For the base case of $k = 1$, it suffices to show that $(\Lambda_1, v_1^{\text{sim}}) \stackrel{d}{=} (\Lambda, v_1^{\text{alg}})$. Note, however, that $v_1^{\text{sim}} = v_1^{\text{alg}} = v_1$, and $\Lambda_1 = (U_1^{\text{sim}})^\top \Lambda (U_1^{\text{sim}}) = U_1(\tilde{v}_1, v_1)^\top \Lambda U_1(\tilde{v}_1, v_1)$. Since v_1 is a deterministic vector, \tilde{v}_1 is independent of Λ , and the distribution of Λ is rotationally invariant, the claim follows.

For the inductive step, assume we know $(\Lambda_k, \{v_j^{\text{sim}}\}_{1 \leq j \leq k}) \stackrel{d}{=} (\Lambda, \{v_j^{\text{alg}}\}_{1 \leq j \leq k})$. Then, note that $v_{k+1}^{\text{alg}} = v_{k+1}(\{\Lambda^i v_j^{\text{alg}}\}_{H_k})$ and $v_{k+1}^{\text{sim}} = v_{k+1}(\{\Lambda_k^i v_j^{\text{sim}}\}_{H_k})$. Thus, we have $v_{k+1}^{\text{alg}} = \mathfrak{D}_{k+1}(\Lambda, \{v_j^{\text{alg}}\}_{1 \leq j \leq k})$ and $v_{k+1}^{\text{sim}} = \mathfrak{D}_{k+1}(\Lambda_k, \{v_j^{\text{sim}}\}_{1 \leq j \leq k})$. In addition, because U_{k+1}^{sim} fixes $\Lambda_k^i v_j^{\text{sim}}$ for all $(i, j) \in H_k$ by (P4), we also have that $\Lambda_{k+1}^i v_j^{\text{sim}} = \Lambda_k^i v_j^{\text{sim}}$ for all $(i, j) \in H_k$, which means $v_{k+1}^{\text{sim}} = \mathfrak{D}_{k+1}(\Lambda_{k+1}, \{v_j^{\text{sim}}\}_{1 \leq j \leq k})$. Therefore, it suffices to show

$$(\Lambda_{k+1}, \{v_j^{\text{sim}}\}_{1 \leq j \leq k}) \stackrel{d}{=} (\Lambda, \{v_j^{\text{alg}}\}_{1 \leq j \leq k}), \quad (5.2)$$

as this implies $(\Lambda_{k+1}, \{v_j^{\text{sim}}\}_{1 \leq j \leq k+1}) \stackrel{d}{=} (\Lambda, \{v_j^{\text{alg}}\}_{1 \leq j \leq k+1})$, which completes the inductive step.

Next, we show that U_{k+1}^{sim} sends \tilde{v}_{k+1} to a random unit vector orthogonal to the simulated queries so far. Note that $\Lambda_{k+1} = (U_{k+1}^{\text{sim}})^\top \Lambda (U_{k+1}^{\text{sim}})$, where, by (P4),

$$U_{k+1}^{\text{sim}} = U_{k+1}(\{\Lambda_k^i v_j^{\text{sim}}\}_{H_k}, (U_{1:k}^{\text{sim}})^\top \tilde{v}_{k+1}, v_{k+1}^{\text{sim}}). \quad (5.3)$$

Note that \tilde{v}_{k+1} has the law of a random unit vector conditional on being orthogonal to $\{\Lambda^i z_j\}_{H_k}$, or equivalently, it is a random unit vector orthogonal to $\{\Lambda^i \tilde{v}_j\}_{H_k}$. Since

$$(U_{1:k}^{\text{sim}})^\top \Lambda^i \tilde{v}_j = (U_{1:k}^{\text{sim}})^\top \Lambda^i (U_{1:k}^{\text{sim}} v_j^{\text{sim}}) = \Lambda_k^i v_j^{\text{sim}}$$

for all $(i, j) \in H_k$ (by (P2)), this means that $(U_{1:k}^{\text{sim}})^\top \tilde{v}_{k+1}$ is orthogonal to $\{\Lambda_k^i v_j^{\text{sim}}\}_{H_k}$. The random direction of \tilde{v}_{k+1} has no dependence on $\{\Lambda^i \tilde{v}_j\}_{H_k}$ apart from being orthogonal to them, which means by (P1), $(U_{1:k}^{\text{sim}})^\top \tilde{v}_{k+1}$ is a *uniformly random* unit vector orthogonal to $\{\Lambda_k^i v_j^{\text{sim}}\}_{H_k}$.

Recalling that $v_{k+1}^{\text{sim}} = \mathfrak{D}_{k+1}(\Lambda_k, \{v_j^{\text{sim}}\}_{1 \leq j \leq k})$, this means that we can rewrite (5.3) as

$$U_{k+1}^{\text{sim}} = U_{k+1}(\{\Lambda_k^i v_j^{\text{sim}}\}_{H_k}, \hat{v}^{\text{sim}}, \mathfrak{D}_{k+1}(\Lambda_k, \{v_j^{\text{sim}}\}_{1 \leq j \leq k})), \quad (5.4)$$

where \hat{v}^{sim} is a random unit vector orthogonal to $\{\Lambda_k^i v_j^{\text{sim}}\}_{H_k}$. As a result, if we define

$$U_{k+1}^{\text{alg}} := U_{k+1}(\{\Lambda^i v_j^{\text{alg}}\}_{H_k}, \hat{v}^{\text{alg}}, \mathfrak{D}_{k+1}(\Lambda, \{v_j^{\text{alg}}\}_{1 \leq j \leq k})), \quad (5.5)$$

where \hat{v}^{alg} is a random unit vector orthogonal to $\{\Lambda^i v_j^{\text{alg}}\}_{H_k}$, then

$$\begin{aligned} (\Lambda_{k+1}, \{v_j^{\text{sim}}\}_{1 \leq j \leq k}) &= ((U_{k+1}^{\text{sim}})^\top \Lambda_k (U_{k+1}^{\text{sim}}), \{v_j^{\text{sim}}\}_{1 \leq j \leq k}) \\ &\stackrel{\text{d}}{=} ((U_{k+1}^{\text{alg}})^\top \Lambda (U_{k+1}^{\text{alg}}), \{v_j^{\text{alg}}\}_{1 \leq j \leq k}). \end{aligned}$$

Above, the first equality follows by definition, and the second follows from our inductive hypothesis that $(\Lambda_k, \{v_j^{\text{sim}}\}_{1 \leq j \leq k}) \stackrel{\text{d}}{=} (\Lambda, \{v_j^{\text{alg}}\}_{1 \leq j \leq k})$, along with (5.4) and (5.5).

We are now in a position to apply the conditioning lemma (Lemma 38). Note that U_{k+1}^{alg} only depends on $\{\Lambda^i v_j^{\text{alg}}\}_{H_k}$ (as well as some randomness in \hat{v}^{alg} , but the randomness is independent of everything else given $\{\Lambda^i v_j^{\text{alg}}\}_{H_k}$, so we can safely condition on it). Hence, we can apply the conditioning lemma with U_{k+1}^{alg} , to obtain that

$$(\Lambda_{k+1}, \{v_j^{\text{sim}}\}_{1 \leq j \leq k}) \stackrel{\text{d}}{=} ((U_{k+1}^{\text{alg}})^\top \Lambda (U_{k+1}^{\text{alg}}), \{v_j^{\text{alg}}\}_{1 \leq j \leq k}) \stackrel{\text{d}}{=} (\Lambda, \{v_j^{\text{alg}}\}_{1 \leq j \leq k}),$$

which establishes (5.2) and thereby concludes the proof. \square

With the block Krylov reduction in hand, we can now establish our second lower bound for sampling from Gaussians.

Theorem 45 (second lower bound for sampling from Gaussians). *There is a universal constant $\epsilon_0 > 0$ such that the query complexity of sampling from Gaussian distributions $\mathcal{N}(0, \Sigma)$ in \mathbb{R}^d , where the condition number κ of Σ satisfies $\kappa \leq d^{1/5-\delta}$, with accuracy ϵ_0 in total variation distance is at least $\Omega_\delta(\sqrt{\kappa} \log d)$.*

Proof. Let U be a random orthogonal matrix, and let $\Lambda = U^\top D U$, $\Lambda' = U^\top D' U$ be as in Lemma 32. We first show that if $\kappa \leq d^{1/5-\delta}$ and c is a sufficiently small constant, no adaptive algorithm that makes less than $c_\delta \sqrt{\kappa} \log d$ queries to the extended oracle can distinguish between Λ and Λ' , with $\Omega(1)$ probability.

First we assume that the algorithm is deterministic, so its behavior is characterized by functions $v_1, v_2(\cdot), \dots, v_K(\cdot)$, as in Lemma 39. The algorithm then proceeds to make queries $v_1^{\text{alg}}, v_2^{\text{alg}}, \dots, v_K^{\text{alg}}$, where $v_k^{\text{alg}} = v_k(\{\Lambda^i v_j^{\text{alg}}\}_{H_{k-1}})$. Lemma 39 shows that the output of the algorithm $\{\Lambda^i v_j^{\text{alg}}\}_{H_K}$ can be entirely simulated by a block Krylov algorithm, which receives $\{\Lambda^i z_k\}_{H_K}$, where z_1, \dots, z_K are

i.i.d. standard Gaussians. Lemma 32 says that a block Krylov algorithm that makes $K = c_\delta \sqrt{\kappa} \log d$ queries, where c_δ is a small constant depending on δ and $\kappa \leq d^{1/5-\delta}$, cannot distinguish between Λ and Λ' with $\Omega(1)$ advantage, which then implies the same for any deterministic algorithm.

If the algorithm is randomized, then it uses a random seed ξ that is independent of Λ and Λ' . So conditional on the random seed, the algorithm will not be able to distinguish Λ and Λ' with $\Omega(1)$ advantage, so the overall probability that the randomized algorithm successfully distinguishes Λ and Λ' also cannot be $\Omega(1)$.

Finally, we note that a sample from $\mathcal{N}(0, \Lambda^{-1})$ versus $\mathcal{N}(0, \Lambda'^{-1})$ can distinguish between the two cases. This means that even if we were able to draw a sample that was $\frac{1}{3}$ -far in total variation distance, we could output the correct answer with probability at least $\frac{2}{3}$. This implies that any sampling algorithm must require at least $\Omega_\delta(\sqrt{\kappa} \log d)$ queries to the extended oracle, and hence at least same number of queries to the standard oracle. \square

Acknowledgments

The authors thank Ainesh Bakshi, Patrik R. Gerber, Piotr Indyk, Thibaut Le Gouic, Philippe Rigollet, Adil Salim, Terence Tao, and Kevin Tian for helpful conversations. SC was supported by the NSF TRIPODS program (award DMS-2022448). JD was supported by a UCLA dissertation year fellowship. CL was supported by the Eric and Wendy Schmidt Center at the Broad Institute of MIT and Harvard. SN was supported by a Google Fellowship, the NSF TRIPODS program (award DMS-2022448), and the NSF Graduate Fellowship.

References

- [AC21] K. Ahn and S. Chewi. “Efficient constrained sampling via the mirror-Langevin algorithm”. In: *Advances in Neural Information Processing Systems*. Ed. by M. Ranzato, A. Beygelzimer, K. Nguyen, P. S. Liang, J. W. Vaughan, and Y. Dauphin. Vol. 34. Curran Associates, Inc., 2021, pp. 28405–28418.
- [AC23] J. M. Altschuler and S. Chewi. “Faster high-accuracy log-concave sampling via algorithmic warm starts”. In: *arXiv preprint 2302.10249* (2023).
- [AT23] J. M. Altschuler and K. Talwar. “Resolving the mixing time of the Langevin algorithm to its stationary distribution for log-concave sampling”. In: *Proceedings of Thirty Sixth Conference on Learning Theory*. Ed. by G. Neu and L. Rosasco. Vol. 195. Proceedings of Machine Learning Research. PMLR, July 2023, pp. 2509–2510.
- [BBH21] M. Brennan, G. Bresler, and B. Huang. “De Finetti-style results for Wishart matrices: combinatorial structure and phase transitions”. In: *arXiv e-prints*, arXiv:2103.14011 (2021).
- [BCESZ22] K. Balasubramanian, S. Chewi, M. A. Erdogdu, A. Salim, and S. Zhang. “Towards a theory of non-log-concave sampling: first-order stationarity guarantees for Langevin Monte Carlo”. In: *Conference on Learning Theory*. PMLR. 2022, pp. 2896–2923.
- [BCW22] A. Bakshi, K. L. Clarkson, and D. P. Woodruff. “Low-rank approximation with $1/\epsilon^{1/3}$ matrix-vector products”. In: *54th Annual ACM SIGACT Symposium on Theory of Computing*. ACM, 2022, pp. 1130–1143.
- [BDER16] S. Bubeck, J. Ding, R. Eldan, and M. Z. Rácz. “Testing for high-dimensional geometry in random graphs”. In: *Random Structures Algorithms* 49.3 (2016), pp. 503–532.

- [Ber18] E. Bernton. “Langevin Monte Carlo and JKO splitting”. In: *Conference on Learning Theory*. PMLR. 2018, pp. 1777–1798.
- [BFG96] Z. Bai, G. Fahey, and G. Golub. “Some large-scale matrix computation problems”. In: *Journal of Computational and Applied Mathematics* 7 (1-2 1996), pp. 71–89.
- [BG18] S. Bubeck and S. Ganguly. “Entropic CLT and phase transition in high-dimensional Wishart matrices”. In: *Int. Math. Res. Not. IMRN* 2 (2018), pp. 588–606.
- [BHSW20] M. Braverman, E. Hazan, M. Simchowitz, and B. E. Woodworth. “The gradient complexity of linear regression”. In: *Conference on Learning Theory, (COLT)*. Vol. 125. Proceedings of Machine Learning Research. PMLR, 2020, pp. 627–647.
- [BKM22] V. Braverman, A. Krishnan, and C. Musco. “Sublinear time spectral density estimation”. In: *54th Annual ACM SIGACT Symposium on Theory of Computing*. ACM, 2022, pp. 1144–1157.
- [Bub15] S. Bubeck. “Convex optimization: algorithms and complexity”. In: *Foundations and Trends® in Machine Learning* 8.3-4 (2015), pp. 231–357.
- [CBL22] N. S. Chatterji, P. L. Bartlett, and P. M. Long. “Oracle lower bounds for stochastic gradient sampling algorithms”. In: *Bernoulli* 28.2 (2022), pp. 1074–1092.
- [CCBJ18] X. Cheng, N. S. Chatterji, P. L. Bartlett, and M. I. Jordan. “Underdamped Langevin MCMC: a non-asymptotic analysis”. In: *Proceedings of the 31st Conference on Learning Theory*. Ed. by S. Bubeck, V. Perchet, and P. Rigollet. Vol. 75. Proceedings of Machine Learning Research. PMLR, July 2018, pp. 300–323.
- [CCSW22] Y. Chen, S. Chewi, A. Salim, and A. Wibisono. “Improved analysis for a proximal algorithm for sampling”. In: *Proceedings of Thirty Fifth Conference on Learning Theory*. Ed. by P.-L. Loh and M. Raginsky. Vol. 178. Proceedings of Machine Learning Research. PMLR, July 2022, pp. 2984–3014.
- [CDWY20] Y. Chen, R. Dwivedi, M. J. Wainwright, and B. Yu. “Fast mixing of Metropolized Hamiltonian Monte Carlo: benefits of multi-step gradients”. In: *J. Mach. Learn. Res.* 21 (2020), pp. 92–1.
- [CE22] Y. Chen and R. Eldan. “Localization schemes: a framework for proving mixing bounds for Markov chains”. In: *arXiv e-prints*, arXiv:2203.04163 (2022).
- [CELSZ22] S. Chewi, M. A. Erdogdu, M. B. Li, R. Shen, and M. Zhang. “Analysis of Langevin Monte Carlo from Poincaré to log-Sobolev”. In: *Proceedings of Thirty Fifth Conference on Learning Theory*. Ed. by P.-L. Loh and M. Raginsky. Vol. 178. Proceedings of Machine Learning Research. PMLR, July 2022, pp. 1–2.
- [CGLGR22] S. Chewi, P. R. Gerber, C. Lu, T. L. Gouic, and P. Rigollet. “The query complexity of sampling from strongly log-concave distributions in one dimension”. In: *Proceedings of Thirty Fifth Conference on Learning Theory*. Ed. by P.-L. Loh and M. Raginsky. Vol. 178. Proceedings of Machine Learning Research. PMLR, July 2022, pp. 2041–2059.
- [CGLL23] S. Chewi, P. R. Gerber, H. Lee, and C. Lu. “Fisher information lower bounds for sampling”. In: *Proceedings of the 34th International Conference on Algorithmic Learning Theory*. Ed. by S. Agrawal and F. Orabona. Vol. 201. Proceedings of Machine Learning Research. PMLR, Feb. 2023, pp. 375–410.
- [Che22] S. Chewi. “Log-concave sampling”. Book draft available at <https://chewisinho.github.io/>. 2022.

- [CKSV18] D. Cohen-Steiner, W. Kong, C. Sohler, and G. Valiant. “Approximating the spectrum of a graph”. In: *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, ACM, 2018, pp. 1263–1271.
- [CLACLR21] S. Chewi, C. Lu, K. Ahn, X. Cheng, T. Le Gouic, and P. Rigollet. “Optimal dimension dependence of the Metropolis-adjusted Langevin algorithm”. In: *Conference on Learning Theory*. PMLR. 2021, pp. 1260–1300.
- [CLLMRS20] S. Chewi, T. Le Gouic, C. Lu, T. Maunu, P. Rigollet, and A. Stromme. “Exponential ergodicity of mirror-Langevin diffusions”. In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 19573–19585.
- [CLW21] Y. Cao, J. Lu, and L. Wang. “Complexity of randomized algorithms for underdamped Langevin dynamics”. In: *Commun. Math. Sci.* 19.7 (2021), pp. 1827–1853.
- [CT06] T. M. Cover and J. A. Thomas. *Elements of information theory*. Second. Wiley-Interscience [John Wiley & Sons], Hoboken, NJ, 2006, pp. xxiv+748.
- [Dal17] A. S. Dalalyan. “Theoretical guarantees for approximate sampling from smooth and log-concave densities”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 79.3 (2017), pp. 651–676.
- [DCWY18] R. Dwivedi, Y. Chen, M. J. Wainwright, and B. Yu. “Log-concave sampling: Metropolis–Hastings algorithms are fast!” In: *Conference on Learning Theory*. PMLR. 2018, pp. 793–797.
- [DK19] A. S. Dalalyan and A. Karagulyan. “User-friendly guarantees for the Langevin Monte Carlo with inaccurate gradient”. In: *Stochastic Processes and their Applications* 129.12 (2019), pp. 5278–5311.
- [DLLW21] Z. Ding, Q. Li, J. Lu, and S. J. Wright. “Random coordinate Langevin Monte Carlo”. In: *Conference on Learning Theory*. PMLR. 2021, pp. 1683–1710.
- [DM17] A. Durmus and E. Moulines. “Nonasymptotic convergence analysis for the unadjusted Langevin algorithm”. In: *The Annals of Applied Probability* 27.3 (2017), pp. 1551–1587.
- [DM21] P. Dharangutte and C. Musco. “Dynamic trace estimation”. In: *Advances in Neural Information Processing Systems* 34. 2021, pp. 30088–30099.
- [DMM19] A. Durmus, S. Majewski, and B. Miasojedow. “Analysis of Langevin Monte Carlo via convex optimization”. In: *J. Mach. Learn. Res.* 20 (2019), Paper No. 73, 46.
- [DR20] A. S. Dalalyan and L. Riou-Durand. “On sampling from a log-concave density using kinetic Langevin diffusions”. In: *Bernoulli* 26.3 (2020), pp. 1956–1988.
- [DT12] A. S. Dalalyan and A. B. Tsybakov. “Sparse regression learning by aggregation and Langevin Monte-Carlo”. In: *J. Comput. System Sci.* 78.5 (2012), pp. 1423–1443.
- [Dvi09] Z. Dvir. “On the size of Kakeya sets in finite fields”. In: *Journal of the American Mathematical Society* 22.4 (2009), pp. 1093–1097.
- [Ede89] A. Edelman. “Eigenvalues and condition numbers of random matrices”. PhD thesis. Department of Mathematics, Massachusetts Institute of Technology, Cambridge, MA, 1989.

- [FYC23] J. Fan, B. Yuan, and Y. Chen. “Improved dimension dependence of a proximal algorithm for sampling”. In: *Proceedings of Thirty Sixth Conference on Learning Theory*. Ed. by G. Neu and L. Rosasco. Vol. 195. Proceedings of Machine Learning Research. PMLR, July 2023, pp. 1473–1521.
- [GLL20] R. Ge, H. Lee, and J. Lu. “Estimating normalizing constants for log-concave distributions: algorithms and lower bounds”. In: *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*. 2020, pp. 579–586.
- [GLL22] S. Gopi, Y. T. Lee, and D. Liu. “Private convex optimization via exponential mechanism”. In: *Proceedings of Thirty Fifth Conference on Learning Theory*. Ed. by P.-L. Loh and M. Raginsky. Vol. 178. Proceedings of Machine Learning Research. PMLR, July 2022, pp. 1948–1989.
- [GV22] K. Gatmiry and S. S. Vempala. “Convergence of the Riemannian Langevin algorithm”. In: *arXiv e-prints*, arXiv:2204.10818 (2022).
- [Hut90] M. F. Hutchinson. “A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines”. In: *Communications in Statistics-Simulation and Computation* 19 (2 1990), pp. 433–450.
- [Jia21] Q. Jiang. “Mirror Langevin Monte Carlo: the case under isoperimetry”. In: *Advances in Neural Information Processing Systems*. Ed. by M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan. Vol. 34. Curran Associates, Inc., 2021, pp. 715–725.
- [JKO98] R. Jordan, D. Kinderlehrer, and F. Otto. “The variational formulation of the Fokker–Planck equation”. In: *SIAM J. Math. Anal.* 29.1 (1998), pp. 1–17.
- [JL15] T. Jiang and D. Li. “Approximation of rectangular beta-Laguerre ensembles and large deviations”. In: *J. Theoret. Probab.* 28.3 (2015), pp. 804–847.
- [Juk11] S. Jukna. *Extremal combinatorics: with applications in computer science*. Vol. 571. Springer, 2011.
- [LST20] Y. T. Lee, R. Shen, and K. Tian. “Logsmooth gradient concentration and tighter runtimes for Metropolized Hamiltonian Monte Carlo”. In: *Conference on Learning Theory*. PMLR. 2020, pp. 2565–2597.
- [LST21a] Y. T. Lee, R. Shen, and K. Tian. “Lower bounds on Metropolized sampling methods for well-conditioned distributions”. In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 18812–18824.
- [LST21b] Y. T. Lee, R. Shen, and K. Tian. “Structured logconcave sampling with a restricted Gaussian oracle”. In: *Conference on Learning Theory*. PMLR. 2021, pp. 2993–3050.
- [LTVW22] R. Li, M. Tao, S. S. Vempala, and A. Wibisono. “The mirror Langevin algorithm converges with vanishing bias”. In: *Proceedings of the 33rd International Conference on Algorithmic Learning Theory*. Ed. by S. Dasgupta and N. Haghtalab. Vol. 167. Proceedings of Machine Learning Research. PMLR, Apr. 2022, pp. 718–742.
- [LV06] L. Lovász and S. Vempala. “Simulated annealing in convex bodies and an $O^*(n^4)$ volume algorithm”. In: *J. Comput. System Sci.* 72.2 (2006), pp. 392–417.
- [MCCFBJ21] Y.-A. Ma, N. S. Chatterji, X. Cheng, N. Flammarion, P. L. Bartlett, and M. I. Jordan. “Is there an analog of Nesterov acceleration for gradient-based MCMC?” In: *Bernoulli* 27.3 (2021), pp. 1942–1992.

- [Mik22] D. Mikulincer. “A CLT in Stein’s distance for generalized Wishart matrices and higher-order tensors”. In: *Int. Math. Res. Not. IMRN* 10 (2022), pp. 7839–7872.
- [MM15] C. Musco and C. Musco. “Randomized block Krylov methods for stronger and faster approximate singular value decomposition”. In: *Advances in Neural Information Processing Systems 28*. 2015, pp. 1396–1404.
- [MMM21] R. A. Meyer, C. Musco, C. Musco, and D. P. Woodruff. “Hutch++: optimal stochastic trace estimation”. In: *4th Symposium on Simplicity in Algorithms*. SIAM, 2021, pp. 142–155.
- [Nes18] Y. Nesterov. *Lectures on convex optimization*. Vol. 137. Springer Optimization and Its Applications. Springer, Cham, 2018, pp. xxiii+589.
- [NS22] A. Nishimura and M. A. Suchard. “Prior-preconditioned conjugate gradient method for accelerated Gibbs sampling in “large n , large p ” Bayesian sparse regression”. In: *Journal of the American Statistical Association* 0.0 (2022), pp. 1–14.
- [NY83] A. S. Nemirovskij and D. B. Yudin. “Problem complexity and method efficiency in optimization”. In: (1983).
- [Per28] O. Perron. “Über einen Satz von Besicovitch”. In: *Mathematische Zeitschrift* 28.1 (1928), pp. 383–386.
- [RC04] C. P. Robert and G. Casella. *Monte Carlo statistical methods*. Second. Springer Texts in Statistics. Springer-Verlag, New York, 2004, pp. xxx+645.
- [RR19] M. Z. Rácz and J. Richey. “A smooth transition from Wishart to GOE”. In: *J. Theoret. Probab.* 32.2 (2019), pp. 898–906.
- [RV08] L. Rademacher and S. Vempala. “Dispersion of mass and the complexity of randomized geometric algorithms”. In: *Adv. Math.* 219.3 (2008), pp. 1037–1069.
- [RWZ20] C. Rashtchian, D. P. Woodruff, and H. Zhu. “Vector-matrix-vector queries for solving linear algebra, statistics, and graph problems”. In: *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*. Vol. 176. LIPIcs. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2020, 26:1–26:20.
- [SAR18] M. Simchowitz, A. E. Alaoui, and B. Recht. “Tight query complexity lower bounds for PCA via finite sample deformed Wigner law”. In: *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*. ACM, 2018, pp. 1249–1259.
- [SL19] R. Shen and Y. T. Lee. “The randomized midpoint method for log-concave sampling”. In: *Advances in Neural Information Processing Systems 32* (2019).
- [SR20] A. Salim and P. Richtarik. “Primal dual interpretation of the proximal stochastic gradient Langevin algorithm”. In: *Advances in Neural Information Processing Systems*. Ed. by H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin. Vol. 33. Curran Associates, Inc., 2020, pp. 3786–3796.
- [SS08] S. Saraf and M. Sudan. “An improved lower bound on the size of Kakeya sets over finite fields”. In: *Analysis & PDE* 1.3 (2008), pp. 375–379.
- [SV14] S. Sachdeva and N. K. Vishnoi. “Faster algorithms via approximation theory”. In: *Found. Trends Theor. Comput. Sci.* 9.2 (2014), pp. 125–210.

- [SWYZ19] X. Sun, D. P. Woodruff, G. Yang, and J. Zhang. “Querying a matrix through matrix-vector products”. In: *46th International Colloquium on Automata, Languages, and Programming*. Vol. 132. LIPIcs. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2019, 94:1–94:16.
- [Sza91] S. J. Szarek. “Condition numbers of random matrices”. In: *J. Complexity* 7.2 (1991), pp. 131–149.
- [Tal19] K. Talwar. “Computational separations between sampling and optimization”. In: *Advances in Neural Information Processing Systems*. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett. Vol. 32. Curran Associates, Inc., 2019.
- [Ver18] R. Vershynin. *High-dimensional probability*. Vol. 47. Cambridge Series in Statistical and Probabilistic Mathematics. An introduction with applications in data science, With a foreword by Sara van de Geer. Cambridge University Press, Cambridge, 2018, pp. xiv+284.
- [VW19] S. Vempala and A. Wibisono. “Rapid convergence of the unadjusted Langevin algorithm: isoperimetry suffices”. In: *Advances in Neural Information Processing Systems 32*. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett. Curran Associates, Inc., 2019, pp. 8094–8106.
- [Wib18] A. Wibisono. “Sampling as optimization in the space of measures: the Langevin dynamics as a composite optimization problem”. In: *Conference on Learning Theory*. PMLR, 2018, pp. 2093–3027.
- [Wib19] A. Wibisono. “Proximal Langevin algorithm: rapid convergence under isoperimetry”. In: *arXiv preprint arXiv:1911.01469* (2019).
- [Woo14] D. P. Woodruff. “Sketching as a tool for numerical linear algebra”. In: *Found. Trends Theor. Comput. Sci.* 10.1-2 (2014), pp. 1–157.
- [WS17] B. Woodworth and N. Srebro. “Lower bound for randomized first order convex optimization”. In: *arXiv e-prints*, arXiv:1709.03594 (2017).
- [WSC22] K. Wu, S. Schmidler, and Y. Chen. “Minimax mixing time of the Metropolis-adjusted Langevin algorithm for log-concave sampling”. In: *Journal of Machine Learning Research* 23.270 (2022), pp. 1–63.
- [WWZ14] K. Wimmer, Y. Wu, and P. Zhang. “Optimal query complexity for estimating the trace of a matrix”. In: *41st International Colloquium on Automata, Languages, and Programming*. Vol. 8572. Lecture Notes in Computer Science. Springer, 2014, pp. 1051–1062.
- [ZPFP20] K. S. Zhang, G. Peyré, J. Fadili, and M. Pereyra. “Wasserstein control of mirror Langevin Monte Carlo”. In: *Proceedings of Thirty Third Conference on Learning Theory*. Ed. by J. Abernethy and S. Agarwal. Vol. 125. Proceedings of Machine Learning Research. PMLR, July 2020, pp. 3814–3841.

A Upper bound for log-concave sampling in constant dimension

In this section we give a simple proof that in constant dimension, one can approximately generate a sample from a log-concave distribution with condition number κ , in $O(\log \kappa)$ queries. Our query dependence also has a polylogarithmic dependence on $\frac{1}{\varepsilon}$, if we wish to generate a sample that is

ε -close in TV distance to the true distribution. (We do not attempt to optimize the dependence on dimension d or the polylogarithmic dependence on $\frac{1}{\varepsilon}$.)

Let V be a convex function that is 1-strongly convex and κ -smooth, such that V is minimized at the origin and $V(0) = 0$. For any real value $y \geq 0$, define $B_V(y)$ to be the set of points x such that $V(x) \leq y$.

First, we note the following basic facts that follow immediately from our convexity assumptions.

Proposition 46 (basic facts about log-concavity). *1. $B_V(y)$ is a convex body for any $y > 0$, and contains 0.*

2. $B_V(y)$ is contained in the ball of radius $\sqrt{2y}$ and contains the ball of radius $\sqrt{2y/\kappa}$.

3. For any $0 < y < y'$, $B_V(y') \subset \frac{y'}{y} B_V(y)$.

Next, we show how to obtain a crude $d^{O(1)}$ -approximation for $B_V(1)$ using $d^{O(1)} \log \kappa$ first-order queries. The proof is essentially folklore and follows from the ellipsoid method.

Proposition 47 (ellipsoid method). *Let B be a convex body that contains $B(0, r)$ and is contained in $B(0, R)$, along with a membership and separation oracle. Using $d^{O(1)} \log \frac{R}{r}$ adaptive queries to the membership and separation oracle, we can find an ellipsoid E centered around some point z such that $E \subset B \subset E'$, where E' is E dilated by an $O(d^{3/2})$ factor about z .*

We can apply the above proposition to the convex body $B_V(1)$.

Corollary 48 (sublevel set approximation). *Using $d^{O(1)} \log \kappa$ adaptive queries to V and ∇V , we can find an ellipsoid E centered around some point z such that $E \subset B_V(1) \subset E'$, where E' is E dilated by an $O(d^{3/2})$ factor about z .*

Proof. It suffices to show that from a single first-order query at a point x , we can generate a membership and separation oracle for $B_V(1)$. Indeed, the membership part is straightforward as we just check whether $V(x) \leq 1$ (which is equivalent to $x \in B_V(1)$). The separation oracle is also simple, and can be done using the gradient. Specifically, suppose that $V(x) > 1$; then, $V(x') \geq V(x) + \langle \nabla V(x), x' - x \rangle$, which means that every x' with $\langle \nabla V(x), x' \rangle \geq \langle \nabla V(x), x \rangle$ is such that $V(x') \geq V(x) > 1$, i.e., $\nabla V(x)$ is a separation oracle for $B_V(1)$ at x . \square

We are able to prove our sampling upper bound, using a rejection sampling approach.

Theorem 49 (upper bound for log-concave sampling). *For any constant $d \geq 2$ and any 1-strongly convex and κ -smooth function V with minimum at 0, we can approximately sample from $\pi \propto \exp(-V)$ to total variation error at most ε using $O(\log \kappa + \log^{O(1)}(1/\varepsilon))$ adaptive queries to V and ∇V (here we emphasize that the asymptotic notation treats d as constant).*

Proof. Given V and any integer $t \geq 1$, let p_t be the probability that a sample from π lies in $tB_V(1)$. The normalizing constant is $Z \geq \int_{B_V(1)} \exp(-V) \geq e^{-1} \text{vol}(B_V(1))$, but integral over $(t+1)B_V(1) \setminus tB_V(1)$ is

$$\int_{(t+1)B_V(1) \setminus tB_V(1)} \exp(-V) \leq \exp(-t) \text{vol}((t+1)B_V(1)) = \exp(-t) (t+1)^d \text{vol}(B_V(1)),$$

using Proposition 46 which implies that $V(x) \geq t$ for any $x \notin tB_V(1)$. Therefore, the probability of $(t+1)B_V(1) \setminus tB_V(1)$ under π is at most

$$\pi((t+1)B_V(1) \setminus tB_V(1)) \leq \frac{\exp(-t) (t+1)^d \text{vol}(B_V(1))}{e^{-1} \text{vol}(B_V(1))} = \exp(-(t-1)) (t+1)^d.$$

By summing this quantity for all integers greater than t , the probability of the complement of $tB_V(1)$ is at most $\sum_{u \geq t} \exp(-(u-1))(u+1)^d = \sum_{u \geq t} \exp(-u + d \log(u+1) + 1)$. Note that for $t \geq \Omega(d \log d)$, this quantity is at most $O(\exp(-t/2))$. Taking $t = C(d \log d + \log(1/\varepsilon))$ for a large constant C , we obtain $\pi(\mathbb{R}^d \setminus tB_V(1)) \leq \varepsilon/2$.

The algorithm now works as follows. We use Corollary 48 to find $E \subset B_V(1) \subset E'$. We pick a uniformly random point X in tE' for $t = C(d \log d + \log(1/\varepsilon))$. We then accept the point X with probability $\exp(-V(X))$, and if we reject we restart the procedure. First, note that this algorithm, upon termination, samples exactly from π conditioned on tE' , which is at most $\frac{\varepsilon}{2}$ away from π in total variation distance. In addition, each rejection sampling step succeeds with probability at least $\text{vol}(E)/(e \text{vol}(tE'))$, since with probability $\text{vol}(E)/\text{vol}(tE')$ we choose a point in E in which case $V(X) \leq 1$ so we accept with probability at least e^{-1} . This is equal to $1/(t O(d^{3/2}))^d = d^{-O(d)} t^{-d} = d^{-O(d)} (\log \frac{1}{\varepsilon})^{-d}$. So, after $(d \log \frac{1}{\varepsilon})^{O(d)}$ rounds of rejection sampling, each of which only needs one query to V , we accept the sample with probability at least $1 - \frac{\varepsilon}{2}$, which means that overall we have generated a sample which is ε -close in distribution to π in total variation distance.

The overall query complexity is a combination of finding E, E' and then running the rejection sampling, for a total complexity of $d^{O(1)} \log \kappa + (d \log \frac{1}{\varepsilon})^{O(d)}$. So, for any fixed dimension d and error probability ε , the query complexity for log-concave sampling is $O(\log \kappa)$. In addition, the dependence on the error probability is polylogarithmic for any fixed d . \square

Remark. We briefly note that the exponential dependence on d is not necessary: using more sophisticated tools developed for sampling from convex bodies one should be able to obtain a complexity of $\log(\kappa) (d \log \frac{1}{\varepsilon})^{O(1)}$. However, we choose to not optimize the dimension dependence in this result for the sake of simplicity, and since we are focused on the setting of $d = O(1)$.

B Upper bound for sampling from Gaussians

Finally, we show a simple proof that, using only $O(\min(\sqrt{\kappa} \log d, d))$ gradient queries, one can generate an approximate sample from a Gaussian $\mathcal{N}(0, \Sigma)$ in d dimensions. Note that the density evaluated at x , up to an additive constant, equals $-\frac{1}{2} x^\top \Lambda x$ for $\Lambda = \Sigma^{-1}$, which means that a gradient query at x amounts to receiving the matrix-vector product Λx .

First, we require a well-known proposition from approximation theory.

Proposition 50 ([SV14, Theorem 3.3]). *For any positive integer s and $0 < \delta < 1$, there exists a polynomial $p_{s,\delta}$ of degree $\lceil \sqrt{2s \ln(2/\delta)} \rceil$ such that $|p_{s,\delta}(x) - x^s| \leq \delta$ for all $x \in [-1, 1]$.*

As a corollary, we have the following result.

Proposition 51 (polynomial approximation of inverse square root). *For any $\kappa \geq 2$ and $\delta < \frac{1}{2}$, there exists a polynomial $q_{\kappa,\delta}$ of degree $O(\sqrt{\kappa} \log \frac{\kappa}{\delta})$ such that $|q_{\kappa,\delta}(x) - x^{-1/2}| \leq \delta/\sqrt{\kappa}$ for all $1 \leq x \leq \kappa$.*

Proof. First, consider the function $(1+x)^{-1/2}$. For $|x| \leq 1 - \frac{1}{\kappa} < 1$, we can use the Taylor series to write

$$(1+x)^{-1/2} = 1 + \sum_{t=1}^{\infty} \frac{(\frac{1}{2}-1)(\frac{1}{2}-2)(\frac{1}{2}-3) \cdots (\frac{1}{2}-t)}{t!} x^t = 1 + \sum_{i=1}^{\infty} c_t x^t,$$

where $|c_t| \leq 1$ for all $t \geq 1$.

Note that for $|x| \leq 1 - \frac{1}{\kappa}$, $|\sum_{t>T} c_t x^t| \leq \sum_{t>T} |x|^t \leq \frac{|x|^T}{1-|x|}$. For $T = O(\kappa \log \frac{\kappa}{\delta})$, we can bound this by $\frac{(1-1/\kappa)^T}{1/\kappa} \leq \frac{\delta}{2}$. Therefore, for all such x ,

$$\left| (1+x)^{-1/2} - \sum_{t=0}^T c_t x^t \right| \leq \frac{\delta}{2},$$

where we have set $c_0 := 1$.

Next, using Proposition 50, we can replace each x^t with $p_{t,\delta}(x)$ where $p_{t,\delta}$ is a polynomial of degree $O(\sqrt{t \log(t/\delta)})$ such that $|p_{t,\delta}(x) - x^t| \leq \delta/(4t^2)$ for all $|x| \leq 1$. (We also let $p_{0,\delta}$ simply be the constant function 1.) Therefore,

$$\left| (1+x)^{-1/2} - \sum_{t=0}^T c_t p_{t,\delta}(x) \right| \leq \frac{\delta}{2} + \sum_{t=1}^T |c_t| \frac{\delta}{4t^2} \leq \delta.$$

In addition, the polynomial $\hat{p} := \sum_{t=0}^T c_t p_{t,\delta}$ has degree at most $O(\sqrt{T \log(T/\delta)}) = O(\sqrt{\kappa \log \frac{\kappa}{\delta}})$.

To finish, $|\hat{p}(x) - x^{-1/2}| \leq \frac{\delta}{\kappa}$ for all $\frac{1}{\kappa} \leq x \leq 1$, which means that

$$\left| \hat{p}\left(\frac{x}{\kappa} - 1\right) \frac{1}{\sqrt{\kappa}} - x^{-1/2} \right| \leq \frac{\delta}{\sqrt{\kappa}} \quad \text{for all } 1 \leq x \leq \kappa.$$

So, there exists a polynomial $q_{\kappa,\delta}$ with $q_{\kappa,\delta}(x) = \hat{p}\left(\frac{x}{\kappa} - 1\right) \frac{1}{\sqrt{\kappa}}$, such that $q_{\kappa,\delta}$ has degree $O(\sqrt{\kappa \log \frac{\kappa}{\delta}})$ and $|q_{\kappa,\delta}(x) - x^{-1/2}| \leq \delta/\sqrt{\kappa}$ for all $1 \leq x \leq \kappa$. \square

We are now ready to prove our query complexity upper bound.

Theorem 52 (optimal algorithm for sampling from Gaussians). *Let $\Lambda = \Sigma^{-1}$ be an unknown positive definite matrix with all eigenvalues between 1 and κ . Then, using $O(\min(\sqrt{\kappa} \log \frac{d}{\varepsilon}, d))$ adaptive matrix-vector queries to Λ , we can produce a sample from a distribution $\hat{\pi}$ such that $\text{KL}(\hat{\pi} \parallel \mathcal{N}(0, \Sigma)) \leq \varepsilon^2$.*

Proof. Choose $X \sim \mathcal{N}(0, I_d)$, define $R = O(\sqrt{\kappa} \log \frac{\kappa}{\delta})$ be the degree of $q_{\kappa,\delta}$, and for simplicity write $q(x) := q_{\kappa,\delta}(x) := \sum_{i=0}^R a_i x^i$. The algorithm works as follows. Using the power method, we compute $X, \Lambda X, \Lambda^2 X, \dots, \Lambda^R X$. We output $Y = \sum_{i=0}^R a_i \Lambda^i X$. Note that $Y \sim \mathcal{N}(0, \hat{\Sigma})$, where we set $\hat{\Sigma} := (\sum_{i=0}^R a_i \Lambda^i)^2$. If $\lambda_1, \dots, \lambda_d$ denote the eigenvalues of Λ , then the eigenvalues of $\hat{\Sigma}$ are $q(\lambda_1)^2, \dots, q(\lambda_d)^2$. The KL divergence is given by

$$\begin{aligned} \text{KL}(\mathcal{N}(0, \hat{\Sigma}) \parallel \mathcal{N}(0, \Sigma)) &\lesssim \sum_{k=1}^d |q(\lambda_k)^2 \lambda_k - 1|^2 \lesssim \sum_{k=1}^d |q(\lambda_k) \lambda_k^{1/2} - 1|^2 \lesssim \sum_{k=1}^d \lambda_k |q(\lambda_k) - \lambda_k^{-1/2}|^2 \\ &\lesssim d\kappa \frac{\delta^2}{\kappa}. \end{aligned}$$

If we set $\delta \asymp \varepsilon/\sqrt{d}$, then we obtain a KL divergence of at most ε^2 .

Finally, we can also learn Λ by querying Λe_i for each unit basis vector e_1, \dots, e_d . So, we can thus learn Σ , and then generate a perfect random sample from $\mathcal{N}(0, \Sigma)$. Hence, the query complexity of generating a sample from $\mathcal{N}(0, \Sigma)$ is at most $O(\min(\sqrt{\kappa} \log \frac{\kappa d}{\varepsilon}, d)) = O(\min(\sqrt{\kappa} \log \frac{d}{\varepsilon}, d))$. \square

Remark. If π is an α -strongly log-concave distribution, then from Pinsker's inequality and Talagrand's transport inequality,

$$\max\{\|\mu - \pi\|_{\text{TV}}^2, \alpha W_2^2(\mu, \pi)\} \lesssim \text{KL}(\mu \parallel \pi).$$

Hence, this algorithmic result for Gaussians complements the two lower bounds in Corollary 19.